

SORT 31 (1) January-June 2007, 3-44

© Institut d'Estadística de Catalunya  
sort@idescat.es

ISSN: 1696-2281

[www.idescat.net/sort](http://www.idescat.net/sort)

# Objective Bayesian point and region estimation in location-scale models

José M. Bernardo

*Universitat de València, Spain*

---

## Abstract

Point and region estimation may both be described as specific *decision problems*. In point estimation, the action space is the set of possible values of the quantity on interest; in region estimation, the action space is the set of its possible credible regions. Foundations dictate that the solution to these decision problems must depend on both the utility function and the prior distribution. Estimators intended for general use should surely be invariant under one-to-one transformations, and this requires the use of an invariant loss function; moreover, an objective solution requires the use of a prior which does not introduce subjective elements. The combined use of an invariant information-theory based loss function, the *intrinsic discrepancy*, and an objective prior, the *reference prior*, produces a general solution to both point and region estimation problems. In this paper, estimation of the two parameters of univariate location-scale models is considered in detail from this point of view, with special attention to the normal model. The solutions found are compared with a range of conventional solutions.

---

MSC: Primary: 62F15, 62C10; secondary: 62F10, 62F15, 62B10

**Keywords:** Confidence Intervals, Credible Regions, Decision Theory, Intrinsic Discrepancy, Intrinsic Loss, Location-Scale Models, Noninformative Prior, Reference Analysis, Region Estimation, Point Estimation.

## 1 Introduction

Point and region estimation of the parameters of location-scale models have a long, fascinating history which is far from settled. Indeed, the list of contributors to the simpler examples of this class of problems, estimation of the normal mean and estimation of the normal variance, reads like a *Who's Who* in 20th century statistics.

---

*Address for correspondence:* José M. Bernardo is Professor of Statistics at the Universitat de València. Departamento de Estadística e I. O., Facultad de Matemáticas, 46100-Burjassot, Valencia, Spain.  
[jose.m.bernardo@uv.es](mailto:jose.m.bernardo@uv.es), [www.uv.es/bernardo](http://www.uv.es/bernardo)  
Received: March 2006

In this paper, an objective Bayesian decision-theoretic solution to both point and region estimation of the parameters of location-scale models is presented, with special attention devoted to the normal model. In marked contrast with most approaches, the solutions found are *invariant* under one-to-one reparametrization.

### 1.1 Notation

Probability distributions are described through their probability density functions, and no notational distinction is made between a random quantity and the particular values that it may take. Bold italic roman fonts are used for observable random vectors (typically data) and bold italic greek fonts for unobservable random vectors (typically parameters); lower case is used for variables and upper case calligraphic for their dominion sets. The standard mathematical convention of referring to functions, say  $f_x(\cdot)$  and  $g_x(\cdot)$  of  $\mathbf{x} \in \mathcal{X}$ , respectively by  $f(\mathbf{x})$  and  $g(\mathbf{x})$  is often used. Thus, the conditional probability density of observable data  $\mathbf{x} \in \mathcal{X}$  given  $\boldsymbol{\omega}$  is represented by either  $p_x(\cdot | \boldsymbol{\omega})$  or  $p(\mathbf{x} | \boldsymbol{\omega})$ , with  $p(\mathbf{x} | \boldsymbol{\omega}) \geq 0$ ,  $\mathbf{x} \in \mathcal{X}$ , and  $\int_{\mathcal{X}} p(\mathbf{x} | \boldsymbol{\omega}) d\mathbf{x} = 1$ , and the posterior density of a non-observable parameter vector  $\boldsymbol{\theta} \in \Theta$  given data  $\mathbf{x}$  is represented by either  $\pi_{\boldsymbol{\theta}}(\cdot | \mathbf{x})$  or  $\pi(\boldsymbol{\theta} | \mathbf{x})$ , with  $\pi(\boldsymbol{\theta} | \mathbf{x}) \geq 0$  and  $\int_{\Theta} \pi(\boldsymbol{\theta} | \mathbf{x}) d\boldsymbol{\theta} = 1$ . Density functions of specific distributions are denoted by appropriate names. In particular, if  $x$  has a normal distribution with mean  $\mu$  and standard deviation  $\sigma$ , its probability density function will be denoted  $N(x | \mu, \sigma)$ , and if  $\lambda$  has a gamma distribution with parameters  $\alpha$  and  $\beta$ , its probability density function will be denoted  $\text{Ga}(\lambda | \alpha, \beta)$ , with  $E[\lambda] = \alpha/\beta$ , and  $\text{Var}[\lambda] = \alpha/\beta^2$ .

It is assumed that available data  $\mathbf{x}$  consist of one observation from the family  $\mathcal{F} \equiv \{p(\mathbf{x} | \boldsymbol{\omega}), \mathbf{x} \in \mathcal{X}, \boldsymbol{\omega} \in \Omega\}$  of probability distributions for  $\mathbf{x} \in \mathcal{X}$ , and that one is interested in point and region estimation of some function  $\boldsymbol{\theta} = \boldsymbol{\theta}(\boldsymbol{\omega}) \in \Theta$  of the unknown parameter vector  $\boldsymbol{\omega}$ . Often, but not necessarily, data consist of a random sample  $\mathbf{x} = \{x_1, \dots, x_n\}$  of some simpler model  $\{q(x | \boldsymbol{\omega}), x \in \mathcal{X}, \boldsymbol{\omega} \in \Omega\}$ , in which case,  $\mathcal{X} = \mathcal{X}^n$  and the likelihood function is  $p(\mathbf{x} | \boldsymbol{\omega}) = \prod_{j=1}^n q(x_j | \boldsymbol{\omega})$ . Without loss of generality, the original parametric family  $\mathcal{F}$  may be written as

$$\mathcal{F} \equiv \{p_x(\cdot | \boldsymbol{\theta}, \boldsymbol{\lambda}), \mathbf{x} \in \mathcal{X}, \boldsymbol{\theta} \in \Theta, \boldsymbol{\lambda} \in \Lambda\} \quad (1)$$

in terms of the vector of interest  $\boldsymbol{\theta}$ , and a vector  $\boldsymbol{\lambda}$  of nuisance parameters. A *point estimator* of  $\boldsymbol{\theta}$  is some function of the data  $\tilde{\boldsymbol{\theta}}(\mathbf{x}) \in \Theta$  such that, for each possible set of observed data  $\mathbf{x}$ ,  $\tilde{\boldsymbol{\theta}}(\mathbf{x})$  could be regarded as an appropriate proxy for the actual, unknown value of  $\boldsymbol{\theta}$ . A *p-credible region* of  $\boldsymbol{\theta}$  is some subset  $C_p(\mathbf{x}, \Theta)$  of  $\Theta$  whose posterior probability is  $p$ . Within this framework, attention in this paper focuses on problems where data consist of a random sample  $\mathbf{x} = \{x_1, \dots, x_n\}$  from a *location-scale* model

$m(x|\mu, \sigma, f)$ , of the form

$$m(x|\mu, \sigma, f) = \sigma^{-1} f\{\sigma^{-1}(x - \mu)\}, \quad x \in \mathfrak{R}, \quad \mu \in \mathfrak{R}, \quad \sigma > 0, \quad (2)$$

where  $f(\cdot)$  is some probability density in  $\mathfrak{R}$ , so that  $f(y) \geq 0$ ,  $\int_{\mathfrak{R}} f(y) dy = 1$ . Interest lies in either the location parameter  $\mu$ , the scale parameter  $\sigma$ , or some one-to-one function of these, and the likelihood function is

$$p(x|\mu, \sigma) = \prod_{j=1}^n m(x_j|\mu, \sigma, f) = \sigma^{-n} \prod_{j=1}^n f\{\sigma^{-1}(x_j - \mu)\}. \quad (3)$$

Standard notation is used for the sample mean and the sample variance, respectively denoted by  $\bar{x} = \sum_{j=1}^n x_j/n$  and  $s^2 = \sum_{j=1}^n (x_j - \bar{x})^2/n$ . Many conventional point estimators of the variance of location-scale models are members of the family of *affine invariant estimators*,

$$\tilde{\sigma}_v^2 = \frac{ns^2}{v} = \frac{1}{v} \sum_{j=1}^n (x_j - \bar{x})^2, \quad v > 0. \quad (4)$$

In particular, with normal data, the MLE of the variance  $\sigma^2$  is  $s^2 = \tilde{\sigma}_n^2$ , and the unbiased estimator is  $\tilde{\sigma}_{n-1}^2$ . More sophisticated estimators may sometimes be defined in terms of affine estimators; for instance, Stein (1964) and Brown (1968) estimators of the normal variance may respectively be written as

$$\tilde{\sigma}_{stein}^2 = \min \left\{ \tilde{\sigma}_{n+1}^2, \tilde{\sigma}_{(n+2)/(1+z^2)}^2 \right\}, \quad \tilde{\sigma}_{brown}^2 = \min \left\{ \tilde{\sigma}_{n-1}^2, \tilde{\sigma}_{n/(1+z^2)}^2 \right\},$$

where  $z = \bar{x}/s$  is the standardized sample mean.

## 1.2 Contents

Section 2 provides a short review of intrinsic estimation, our approach to both point and region estimation. An information-theory based invariant loss function, the *intrinsic discrepancy*, is proposed as a reasonable general alternative to the conventional (non-invariant) quadratic loss. As is usually the case in modern literature, point estimation is described as a decision problem where the action space is the set of possible values for the quantity of interest; an intrinsic point estimator is then defined as the Bayes estimator which corresponds to the intrinsic loss and the appropriate reference prior. This provides a general *invariant* objective Bayes point estimator. Less conventionally, region estimation is also described as a decision problem where, for each  $p$ , the action space is the set of possible  $p$ -credible regions for the quantity of interest; a  $p$ -credible intrinsic region estimator is then defined as the lowest posterior loss  $p$ -credible region with respect to the intrinsic loss and the appropriate reference prior. This provides a

general *invariant* objective Bayes region estimator which always contains the intrinsic point estimator.

In Section 3 location-scale models are analyzed from this point of view. In particular, intrinsic point estimators and intrinsic region estimators are derived for the mean of a normal model, the variance of a normal model, and the scale parameter of a Cauchy model.

## 2 Intrinsic Estimation

### 2.1 The intrinsic discrepancy loss function

Point estimation of some parameter vector  $\theta \in \Theta$  is customarily described as a decision problem where the action space is the set  $\mathcal{A} = \{\tilde{\theta}; \tilde{\theta} \in \Theta\}$  of possible values of the vector of interest. Foundations dictate (see e.g., Bernardo and Smith, 1994, Ch. 2 and references therein) that to solve this decision problem it is necessary to specify a *loss function*  $\ell\{\tilde{\theta}, \theta\}$ , such that  $\ell\{\tilde{\theta}, \theta\} \geq 0$  and  $\ell\{\theta, \theta\} = 0$ , which describes, as a function of  $\theta$ , the loss suffered from using  $\tilde{\theta}$  as a proxy for the unknown value of  $\theta$ . The loss function is context specific, and should be chosen in terms of the anticipated uses of the estimate; however, a number of conventional loss functions have been suggested for those situations where no particular uses are envisaged, as in scientific inference. The simplest of these conventional loss functions (which typically ignore the presence nuisance parameters) is the ubiquitous *quadratic loss*,  $\ell\{\tilde{\theta}, (\theta, \lambda)\} = (\tilde{\theta} - \theta)'(\tilde{\theta} - \theta)$ ; the corresponding Bayes estimator, if this exists, is the *posterior mean*,  $E[\theta | \mathbf{x}]$ . Another common conventional loss function is the *zero-one loss*, defined as  $\ell\{\tilde{\theta}, (\theta, \lambda)\} = 1$ , if  $\tilde{\theta}$  does not belong to a  $\epsilon$ -radius neighbourhood of  $\theta$ , and zero otherwise; as  $\epsilon \rightarrow 0$ , the corresponding Bayes estimator converges to the posterior mode,  $\text{Mo}[\theta | \mathbf{x}]$ . For details, see, e.g., Bernardo and Smith (1994, p. 257).

**Example 1 (Normal variance)** With the usual objective prior  $\pi(\mu, \sigma) = \sigma^{-1}$ , the (marginal) reference posterior density of  $\sigma$  is the square root inverted gamma

$$\pi(\sigma | \mathbf{x}) = \pi(\sigma | s, n) = \frac{n^{(n-1)/2} s^{n-1}}{2^{(n-3)/2} \Gamma[(n-1)/2]} \sigma^{-n} e^{-\frac{1}{2}n s^2 / \sigma^2}, \quad n \geq 2. \quad (5)$$

The quadratic loss in terms of the variance,  $\ell\{\tilde{\sigma}^2, \sigma^2\} = (\tilde{\sigma}^2 - \sigma^2)^2$ , leads to  $E[\sigma^2 | \mathbf{x}] = \tilde{\sigma}_{n-3}^2$  (which obviously requires  $n \geq 3$ ). Similarly, the quadratic loss in terms of the standard deviation,  $\ell\{\tilde{\sigma}, \sigma\} = (\tilde{\sigma} - \sigma)^2$ , yields

$$E[\sigma | \mathbf{x}] = \sqrt{\frac{n}{2}} \frac{\Gamma[(n-2)/2]}{\Gamma[(n-1)/2]} s, \quad n \geq 3. \quad (6)$$

For moderate  $n$  values, the Stirling approximation of the Gamma functions in (6) produces  $E[\sigma | \mathbf{x}]^2 \approx \tilde{\sigma}_{n-5/2}^2$ . Notice that, using the conventional quadratic loss function, the Bayes estimate of  $\sigma^2$  is *not* the square of the Bayes estimate of  $\sigma$ . This lack of invariance is not an special feature of the quadratic loss; on the contrary, this is the case of most conventional loss functions. For instance, the use of the slightly more sophisticated standardized quadratic loss function on the variance,

$$\ell_{stq}(\tilde{\sigma}^2, \sigma^2) = [(\tilde{\sigma}^2/\sigma^2) - 1]^2 \quad (7)$$

yields (if  $n \geq 2$ , for  $\pi(\sigma | \mathbf{x})$  to be proper)

$$\arg \min_{\tilde{\sigma}^2 > 0} \int_0^\infty [(\tilde{\sigma}^2/\sigma^2) - 1]^2 \pi(\sigma | \mathbf{x}) d\sigma = \frac{n s^2}{n+1} = \tilde{\sigma}_{n+1}^2, \quad (8)$$

which is also the minimum risk equivariant estimator (MRIE) of  $\sigma^2$  under this loss, while the standardized quadratic loss function in terms of the standard deviation,  $\ell(\tilde{\sigma}, \sigma) = [(\tilde{\sigma}/\sigma) - 1]^2$  yields (again, if  $n \geq 2$ )

$$\arg \min_{\tilde{\sigma} > 0} \int_0^\infty [(\tilde{\sigma}/\sigma) - 1]^2 \pi(\sigma | \mathbf{x}) d\sigma = \sqrt{\frac{n}{2}} \frac{\Gamma[n/2]}{\Gamma[(n+1)/2]} s, \quad (9)$$

which is different from (6), and whose square is *not* (8). Similarly, for the zero-one loss in terms of  $\sigma^2$ , the Bayes estimator is the mode of the posterior distribution of  $\sigma^2$ ,  $\pi(\sigma^2 | \mathbf{x}) = \pi(\sigma | \mathbf{x})/(2\sigma)$ , which is  $\text{Mo}(\sigma^2 | \mathbf{x}) = \tilde{\sigma}_{n+1}^2$ , the same as (8), while the Bayes estimator for the zero-one loss in terms of  $\sigma$  is  $\text{Mo}(\sigma | \mathbf{x}) = s$ , the MLE of  $\sigma$ , whose square is obviously not the same as (8). For further information on alternative point estimators of the normal variance, see Brewster and Zidek (1974) and Rukhin (1987).

As Example 1 dramatically illustrates, conventional loss functions are typically *not* invariant under reparametrization. As a consequence, the Bayes estimator  $\phi^*$  of a one-to-one transformation  $\phi = \phi(\theta)$  of the original parameter  $\theta$  is not necessarily  $\phi(\theta^*)$  and thus, for each loss function, one may produce as many *different* estimators of the same quantity as alternative parametrizations one is prepared to consider, a less than satisfactory situation. Indeed, scientific applications *require* this type of invariance. It would certainly be hard to argue that the best estimate of, say the age of the universe is  $\theta^*$  but that the best estimate of the logarithm of that age is *not*  $\log(\theta^*)$ . Invariant loss functions are required to guarantee invariant estimators.

With no nuisance parameters, *intrinsic loss functions* (Robert, 1996), of the general form  $\ell(\tilde{\theta}, \theta) = \ell\{p_x(\cdot | \tilde{\theta}), p_x(\cdot | \theta)\}$  shift attention from the discrepancy between the estimate  $\tilde{\theta}$  and the true value  $\theta$ , to the more relevant discrepancy between the statistical *models* they label, and they are always invariant under one-to-one reparametrization. The *intrinsic discrepancy*, introduced by Bernardo and Rueda (2002), is a particular intrinsic loss with specially attractive properties.

**Definition 1 (Intrinsic Discrepancy)** The intrinsic discrepancy between two elements  $p_x(\cdot | \omega_1)$  and  $p_x(\cdot | \omega_2)$  of the parametric family of distributions  $\mathcal{F} = \{p_x(\cdot | \omega), \mathbf{x} \in \mathcal{X}(\omega), \omega \in \Omega\}$ , is

$$\begin{aligned}\delta_x(\omega_1, \omega_2) &= \delta\{p_x(\cdot | \omega_1), p_x(\cdot | \omega_2)\} = \min\{\kappa_x(\omega_1 | \omega_2), \kappa_x(\omega_2 | \omega_1)\}, \\ \kappa_x(\omega_j | \omega_i) &= \int_{\mathcal{X}(\omega_i)} p_x(\mathbf{x} | \omega_i) \log \frac{p_x(\mathbf{x} | \omega_i)}{p_x(\mathbf{x} | \omega_j)} d\mathbf{x},\end{aligned}$$

The intrinsic discrepancy  $\delta_x\{\mathcal{F}_1, \mathcal{F}_2\}$  between two subsets  $\mathcal{F}_1$  and  $\mathcal{F}_2$  of  $\mathcal{F}$  is the minimum intrinsic discrepancy between its elements,

$$\delta_x(\mathcal{F}_1, \mathcal{F}_2) = \min_{\omega_1 \in \mathcal{F}_1, \omega_2 \in \mathcal{F}_2} \delta\{p_x(\cdot | \omega_1), p_x(\cdot | \omega_2)\}$$

Thus, the intrinsic discrepancy  $\delta(\omega_1, \omega_2)$  between two parameter values  $\omega_1$  and  $\omega_2$  is the minimum Kullback-Leibler directed logarithmic divergence (Kullback and Leibler, 1951) between the distributions  $p_x(\cdot | \omega_1)$  and  $p_x(\cdot | \omega_2)$  which they label. Notice that this is obviously independent of the particular parametrization chosen to describe the distributions. The intrinsic discrepancy is a divergence measure in the class  $\mathcal{F}$ ; indeed, (i) it is symmetric, (ii) it is non-negative and (iii) it is zero if, and only if,  $p_x(\mathbf{x} | \omega_1) = p_x(\mathbf{x} | \omega_2)$  almost everywhere. Notice that in Definition 1 the possible dependence of the sampling space  $\mathcal{X} = \mathcal{X}(\omega)$  on the parameter value  $\omega$  is explicitly allowed, so that the intrinsic discrepancy may be used with non-regular models where the support  $\mathcal{X}(\omega_1)$  of, say,  $p_x(\cdot | \omega_1)$  may be strictly smaller than the support  $\mathcal{X}(\omega_2)$  of  $p_x(\cdot | \omega_2)$ .

The intrinsic discrepancy is also invariant under one-to-one transformations of the random vector  $\mathbf{x}$ . Moreover, directed logarithmic divergences are *additive* with respect to conditionally independent observations. Consequently, if  $\mathbf{x} = \{x_1, \dots, x_n\}$  is a random sample from, say  $q_x(\cdot | \omega)$  so that the probability model is  $p(\mathbf{x} | \omega) = \prod_{i=1}^n q(x_i | \omega)$ , then the intrinsic discrepancy  $\delta_x\{\omega_1, \omega_2\}$  between  $p_x(\cdot | \omega_1)$  and  $p_x(\cdot | \omega_2)$  is simply  $n \delta_x\{\omega_1, \omega_2\}$ , that is,  $n$  times the intrinsic discrepancy between  $q_x(\cdot | \omega_1)$  and  $q_x(\cdot | \omega_2)$ .

In the context of point estimation, the intrinsic discrepancy leads naturally to the (invariant) intrinsic discrepancy loss  $\delta_x\{\tilde{\theta}, (\theta, \lambda)\}$  defined as the intrinsic discrepancy between the assumed model  $p_x(\cdot | \theta, \lambda)$  and its closest approximation within the set  $\{p_x(\cdot | \tilde{\theta}, \tilde{\lambda}), \tilde{\lambda} \in \Lambda\}$  of all models with  $\theta = \tilde{\theta}$ .

**Definition 2 (Intrinsic discrepancy loss)** Consider the family of probability distributions  $\mathcal{F} = \{p_x(\cdot | \theta, \lambda), \theta \in \Theta, \lambda \in \Lambda, \mathbf{x} \in \mathcal{X}(\omega, \lambda)\}$ . The intrinsic discrepancy loss from using  $\tilde{\theta}$  as a proxy for  $\theta$  is

$$\delta_x\{\tilde{\theta}, (\theta, \lambda)\} = \inf_{\tilde{\lambda} \in \Lambda} \delta_x\{(\tilde{\theta}, \tilde{\lambda}), (\theta, \lambda)\},$$

the intrinsic discrepancy between  $p_x(\cdot | \theta, \lambda)$  and the set  $\{p_x(\cdot | \tilde{\theta}, \tilde{\lambda}), \tilde{\lambda} \in \Lambda\}$ .

Notice that the value of  $\delta_x\{\tilde{\theta}, (\theta, \lambda)\}$  does *not* depend on the particular parametrization chosen to describe the problem. Indeed, for any one-to-one reparametrizations  $\phi = \phi(\theta)$  and  $\psi = \psi(\lambda)$ ,

$$\delta_x\{\tilde{\phi}, (\phi, \psi)\} = \delta_x\{\tilde{\theta}, (\theta, \lambda)\} \quad (10)$$

so that, as one should surely require, the loss suffered from using  $\tilde{\phi} = \phi(\tilde{\theta})$  as a proxy for  $\phi(\theta)$  is precisely the same as the loss suffered from using  $\tilde{\theta}$  as a proxy for  $\theta$ , and this is true for any parametrization of the nuisance parameter vector.

Under frequently met regularity conditions, the two minimizations required in Definition 2 may be interchanged. This makes analytical derivation of the intrinsic loss considerably simpler.

**Theorem 1** (*Computation of the intrinsic discrepancy loss*) *Let  $\mathcal{F}$  be a parametric family of probability distributions*

$$\mathcal{F} = \{p(x | \theta, \lambda), \theta \in \Theta, \lambda \in \Lambda, x \in \mathcal{X}(\theta, \lambda)\},$$

*with convex support  $\mathcal{X}(\theta, \lambda)$  for all  $\theta$  and  $\lambda$ . Then,*

$$\begin{aligned} \delta_x\{\tilde{\theta}, (\theta, \lambda)\} &= \inf_{\tilde{\lambda} \in \Lambda} \min \left\{ \kappa_x\{\tilde{\theta}, \tilde{\lambda} | \theta, \lambda\}, \kappa_x\{\theta, \lambda | \tilde{\theta}, \tilde{\lambda}\} \right\} \\ &= \min \left\{ \inf_{\tilde{\lambda} \in \Lambda} \kappa_x\{\tilde{\theta}, \tilde{\lambda} | \theta, \lambda\}, \inf_{\tilde{\lambda} \in \Lambda} \kappa_x\{\theta, \lambda | \tilde{\theta}, \tilde{\lambda}\} \right\} \end{aligned}$$

*Proof.* This follows from the fact that, if  $\mathcal{X}(\theta, \lambda)$  is a convex set, then the two directed logarithmic divergences involved in the definition are convex functionals. For details, see Juárez (2004).  $\square$

**Example 2** (*Normal variance, continued*) Consider  $\mathbf{x} = \{x_1, \dots, x_n\}$ . The directed logarithmic divergence of  $\prod_{i=1}^n N(x_i | \tilde{\mu}, \tilde{\sigma})$  from  $\prod_{i=1}^n N(x_i | \mu, \sigma)$  is

$$\begin{aligned} \kappa_x\{\tilde{\mu}, \tilde{\sigma} | \mu, \sigma\} &= n \int_{-\infty}^{\infty} N(x | \mu, \sigma) \log \left[ \frac{N(x | \mu, \sigma)}{N(x | \tilde{\mu}, \tilde{\sigma})} \right] dx \\ &= \frac{n}{2} \left[ \frac{\sigma^2}{\tilde{\sigma}^2} - 1 - \log \frac{\sigma^2}{\tilde{\sigma}^2} + \frac{(\mu - \tilde{\mu})^2}{\tilde{\sigma}^2} \right]. \end{aligned} \quad (11)$$

This is minimized when  $\tilde{\mu} = \mu$ , to yield

$$\inf_{\tilde{\mu} \in \mathbb{R}} \kappa_x\{\tilde{\mu}, \tilde{\sigma} | \mu, \sigma\} = \frac{n}{2} g \left( \frac{\sigma^2}{\tilde{\sigma}^2} \right) = \frac{n}{2} g(\phi),$$

where  $\phi = \tilde{\sigma}^2/\sigma^2$ , and  $g(\cdot)$  is the *linlog* function defined by

$$g(t) = (t - 1) - \log t, \quad t > 0. \quad (12)$$

Notice that  $g(t) \geq 0$  and  $g(t) = 0$  if (and only if)  $t = 1$ . This follows from the fact that  $g(t)$  is the absolute distance between  $\log t$  and its tangent at  $t = 1$ .

Exchanging the roles of  $(\tilde{\mu}, \tilde{\sigma})$  and  $(\mu, \sigma)$ , it is similarly found that  $\kappa_x\{\mu, \sigma | \tilde{\mu}, \tilde{\sigma}\}$  is also minimized when  $\tilde{\mu} = \mu$ , to yield

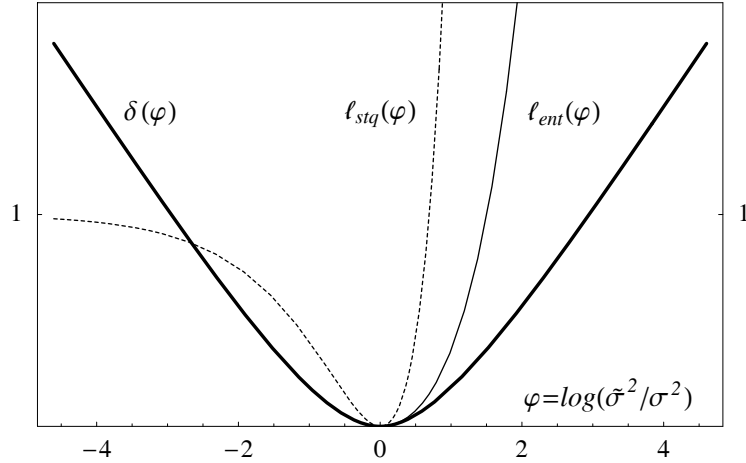
$$\inf_{\tilde{\mu} \in \mathbb{R}} \kappa_x\{\mu, \sigma | \tilde{\mu}, \tilde{\sigma}\} = \frac{n}{2} g\left(\frac{\tilde{\sigma}^2}{\sigma^2}\right) = \frac{n}{2} g\left(\frac{1}{\phi}\right).$$

Moreover,  $g(t) < g(1/t)$  if, and only if,  $t < 1$  and hence, using Theorem 1,

$$\delta_x\{\tilde{\sigma}, (\mu, \sigma)\} = \delta_x\{\phi\} = \begin{cases} \frac{n}{2} g(\phi) & \text{if } \phi < 1, \\ \frac{n}{2} g(1/\phi) & \text{if } \phi \geq 1, \end{cases} \quad \phi = \frac{\tilde{\sigma}^2}{\sigma^2}. \quad (13)$$

Thus, for fixed  $n$ , the intrinsic discrepancy loss  $\delta_x\{\tilde{\sigma}, (\mu, \sigma)\}$  only depends on the ratio  $\phi = \tilde{\sigma}^2/\sigma^2$ . The intrinsic discrepancy loss is closely related to the (also invariant) *entropy* loss,

$$\ell_{ent}\{\tilde{\sigma}, \sigma\} = \ell_{ent}\{\phi\} = \int_{-\infty}^{\infty} N(x | \mu, \sigma) \log \left[ \frac{N(x | \mu, \sigma)}{N(x | \mu, \tilde{\sigma})} \right] dx = \frac{1}{2} g(\phi), \quad (14)$$



**Figure 1:** Intrinsic discrepancy loss (solid line), entropy loss (continuous line), and standardized quadratic loss (dotted line) for point estimation of the normal variance, as a function of  $\psi = \log(\tilde{\sigma}^2/\sigma^2)$ .



which Brown (1990) attributes to Stein. Except for the proportionality constant  $n$  (which does not affect estimation), the entropy loss (14) is the same as the intrinsic discrepancy loss (13) whenever  $\tilde{\sigma} < \sigma$ . Indeed, the intrinsic discrepancy loss may be seen as a symmetrized version of the entropy loss.

Notice that, for all values of the ratio  $\phi = \tilde{\sigma}^2/\sigma^2$ ,  $\delta_x\{\phi\} = \delta_x\{1/\phi\}$ ; hence, the intrinsic loss *equally* penalizes overestimation and underestimation. In sharp contrast, both the entropy loss and the often recommended standardized quadratic loss function, which is also a function of the ratio  $\phi$ ,

$$\ell_{sq}(\tilde{\sigma}^2, \sigma^2) = [(\tilde{\sigma}^2/\sigma^2) - 1]^2 = (\phi - 1)^2,$$

clearly underpenalize small estimators, thus yielding estimators of the variance which are too small. This is illustrated in Figure 1, where the functions  $\delta_x(\phi)$  (for  $n = 1$ ),  $\ell_{ent}\{\phi\}$ , and  $\ell_{sq}(\phi)$  are all represented as a function of  $\varphi = \log \phi$ . More conventional loss functions, as the usual quadratic loss,

$$\ell_{quad}(\tilde{\sigma}^2, \sigma^2) = [\tilde{\sigma}^2 - \sigma^2]^2 = \sigma^4(\phi - 1)^2,$$

are not even invariant with respect to affine transformations. All this led Stein (1964, p. 156) to write “I find it hard to take the problem of estimating  $\sigma^2$  with quadratic loss very seriously”.

## 2.2 Reference posterior expectation

Given data  $\mathbf{x}$  generated by  $p(\mathbf{x}|\boldsymbol{\theta}, \boldsymbol{\lambda})$ , a situation with no prior information about the value of  $\boldsymbol{\theta}$  is formally described by the *reference prior*  $\pi(\boldsymbol{\lambda}|\boldsymbol{\theta})\pi(\boldsymbol{\theta})$  which corresponds to the model  $p_x(\cdot|\boldsymbol{\theta}, \boldsymbol{\lambda})$  when  $\boldsymbol{\theta}$  is the quantity of interest (Bernardo, 1979; Berger and Bernardo, 1992; Bernardo, 2005a). In this case, all available information about  $\boldsymbol{\theta}$  is encapsulated in its (marginal) reference posterior distribution,  $\pi(\boldsymbol{\theta}|\mathbf{x}) = \int_{\Lambda} \pi(\boldsymbol{\theta}, \boldsymbol{\lambda}|\mathbf{x}) d\boldsymbol{\lambda}$  where, by Bayes theorem,  $\pi(\boldsymbol{\theta}, \boldsymbol{\lambda}|\mathbf{x}) \propto p(\mathbf{x}|\boldsymbol{\theta}, \boldsymbol{\lambda})\pi(\boldsymbol{\lambda}|\boldsymbol{\theta})\pi(\boldsymbol{\theta})$ . If numerical summaries of the information encapsulated in  $\pi(\boldsymbol{\theta}|\mathbf{x})$  are further required in the form of either point or region estimators of  $\boldsymbol{\theta}$  under some specified loss function  $\ell\{\tilde{\boldsymbol{\theta}}, (\boldsymbol{\theta}, \boldsymbol{\lambda})\}$ , then the reference posterior expected loss

$$l(\tilde{\boldsymbol{\theta}}|\mathbf{x}) = \int_{\Theta} \int_{\Omega} \ell\{\tilde{\boldsymbol{\theta}}, (\boldsymbol{\theta}, \boldsymbol{\lambda})\} \pi(\boldsymbol{\theta}, \boldsymbol{\lambda}|\mathbf{x}) d\boldsymbol{\theta} d\boldsymbol{\lambda}$$

(from using  $\tilde{\boldsymbol{\theta}}$  as a proxy for  $\boldsymbol{\theta}$ ) has to be evaluated. In view of the arguments given above, attention will focus on the intrinsic discrepancy reference expected loss, or *intrinsic expected loss*, for short.

**Definition 3 (Intrinsic expected loss)** Consider the parametric family of probability distributions

$$\mathcal{F} = \{p_x(\cdot | \boldsymbol{\theta}, \boldsymbol{\lambda}), \mathbf{x} \in \mathcal{X}(\boldsymbol{\theta}, \boldsymbol{\lambda}), \boldsymbol{\theta} \in \boldsymbol{\Theta}, \boldsymbol{\lambda} \in \boldsymbol{\Lambda}, \}, \quad (15)$$

The intrinsic expected loss from using  $\tilde{\boldsymbol{\theta}}$  given data  $\mathbf{x}$ , denoted  $d(\tilde{\boldsymbol{\theta}} | \mathbf{x})$ , is the posterior expectation of the intrinsic discrepancy loss,  $\delta_x\{\tilde{\boldsymbol{\theta}}, (\boldsymbol{\theta}, \boldsymbol{\lambda})\}$  with respect to the joint reference posterior,  $\pi(\boldsymbol{\theta}, \boldsymbol{\lambda} | \mathbf{x})$ ,

$$d(\tilde{\boldsymbol{\theta}} | \mathbf{x}) = \int_{\boldsymbol{\Theta}} \int_{\boldsymbol{\Lambda}} \delta_x\{\tilde{\boldsymbol{\theta}}, (\boldsymbol{\theta}, \boldsymbol{\lambda})\} \pi(\boldsymbol{\theta}, \boldsymbol{\lambda} | \mathbf{x}) d\boldsymbol{\theta} d\boldsymbol{\lambda}, \quad (16)$$

where  $\pi(\boldsymbol{\theta}, \boldsymbol{\lambda} | \mathbf{x}) \propto p(\mathbf{x} | \boldsymbol{\theta}, \boldsymbol{\omega}) \pi(\boldsymbol{\lambda} | \boldsymbol{\theta}) \pi(\boldsymbol{\theta})$ , and  $\pi(\boldsymbol{\lambda} | \boldsymbol{\theta}) \pi(\boldsymbol{\theta})$  is the (joint) reference prior when  $\boldsymbol{\theta}$  is the quantity of interest.

The function  $d(\tilde{\boldsymbol{\theta}} | \mathbf{x})$  measures the posterior expected loss from using  $\tilde{\boldsymbol{\theta}}$  as a proxy for the unknown value of  $\boldsymbol{\theta}$ , in terms of the expected intrinsic discrepancy between the assumed model,  $p_x(\cdot | \boldsymbol{\theta}, \boldsymbol{\lambda})$ , and the class

$$\mathcal{F}_{\tilde{\boldsymbol{\theta}}} = \{p(\mathbf{x} | \tilde{\boldsymbol{\theta}}, \tilde{\boldsymbol{\lambda}}), \tilde{\boldsymbol{\lambda}} \in \boldsymbol{\Lambda}, \mathbf{x} \in \mathcal{X}(\boldsymbol{\omega}, \boldsymbol{\lambda}) \quad (17)$$

of those models in  $\mathcal{F}$  for which  $\boldsymbol{\theta} = \tilde{\boldsymbol{\theta}}$ .

The intrinsic expected loss provides an objective measure of the *compatibility* of the value  $\tilde{\boldsymbol{\theta}}$  with the observed data  $\mathbf{x}$ , with a nice interpretation in terms of likelihood ratios. Indeed, it immediately follows from Definitions 1, 2 and 3, that  $d(\tilde{\boldsymbol{\theta}} | \mathbf{x})$  is the posterior expectation of the minimum expected log-likelihood ratio between the true model and the closest model for which  $\boldsymbol{\theta} = \tilde{\boldsymbol{\theta}}$ . For instance, if  $d(\tilde{\boldsymbol{\theta}} | \mathbf{x}) = \log 100$ , then data  $\mathbf{x}$  are expected to be about 100 times more likely under the true (unknown) model than under any model within this family with  $\boldsymbol{\theta} = \tilde{\boldsymbol{\theta}}$ .

**Example 3 (Normal variance, continued)** In Example 2 the intrinsic discrepancy loss from using  $\tilde{\sigma}$  as a proxy for  $\sigma$  was seen to be a function  $\delta_x(\phi)$  of the ratio  $\phi = \tilde{\sigma}^2 / \sigma^2$  (Equation 13). Changing variables in (5), the reference posterior of  $\phi$  is

$$\pi(\phi | \mathbf{x}) = \text{Ga}(\phi | \frac{n-1}{2}, \frac{n s^2}{2 \tilde{\sigma}^2}). \quad (18)$$

Hence, the intrinsic expected loss from using  $\tilde{\sigma}$  as a proxy for  $\sigma$  is

$$d(\tilde{\sigma} | \mathbf{x}) = d(\tilde{\sigma} | s^2, n) = \int_0^\infty \delta_x(\phi) \text{Ga}(\phi | \frac{n-1}{2}, \frac{n s^2}{2 \tilde{\sigma}^2}) d\phi \quad (19)$$

which may easily be computed by one-dimensional numerical integration. Good analytical approximations will however be provided in Section 3. As a numerical

illustration, a random sample of size  $n = 12$  was simulated from a normal distribution with  $\mu = 5$  and  $\sigma = 2$ , yielding  $\bar{x} = 4.214$  and  $s = 2.071$ . The corresponding intrinsic expected loss  $d(\tilde{\sigma}|\mathbf{x})$ , represented in the lower panel of Figure 2, is locally convex around a unique minimum.

### 2.3 Intrinsic Point and Region Estimation

Bayes estimates are, by definition, those which minimize the expected posterior loss. The *intrinsic estimate* is the Bayes estimate which corresponds to the intrinsic discrepancy loss and the reference posterior distribution, i.e., that value  $\tilde{\theta}_{int}(\mathbf{x}) \in \Theta$  which minimizes the intrinsic expected loss. Formally,

**Definition 4 (Intrinsic point estimator)** Consider again the parametric family of probability distributions  $\mathcal{F}$  defined by (15). An intrinsic estimator of  $\theta$  is a value

$$\tilde{\theta}_{int}(\mathbf{x}) = \min_{\tilde{\theta} \in \Theta} d\{\tilde{\theta}|\mathbf{x}\},$$

which minimizes the intrinsic discrepancy reference posterior loss (16).

Under general regularity conditions, the intrinsic expected loss  $d\{\tilde{\theta}|\mathbf{x}\}$  is locally convex near its absolute minimum and, therefore, the intrinsic estimate typically exists and it is unique. Moreover, since both the intrinsic loss function and the reference prior are invariant under one-to-one reparametrization, the intrinsic estimator  $\tilde{\psi}_{int}(\mathbf{x})$  of any one-to-one function  $\psi(\theta)$  of  $\theta$  will simply be  $\tilde{\psi}_{int} = \psi(\tilde{\theta}_{int})$ . For more details on intrinsic estimation, see Bernardo and Juárez (2003).

Bayesian region estimation is typically based on posterior credible regions, i.e., sets of  $\theta$  values with pre-specified posterior probabilities. However, for any fixed  $p$ , there are typically infinitely many  $p$ -credible regions. In most cases, these are chosen to be either highest posterior density (HPD) regions, or probability centred regions.

It is well known that the ubiquitous *highest posterior density* (HPD) regions are *not* consistent under reparametrization. Thus, if  $\phi\{\theta\}$  is a one-to-one function of  $\theta$ , the image of a HPD  $p$ -credible region of  $\theta$  will *not* generally be HPD for  $\phi$ . Thus, if  $C_p$  is a HPD  $p$ -credible set estimate for, say, the perihelion of the Earth,  $\log(C_p)$  will *not* be a HPD  $p$ -credible set estimate for its logarithm. This suggests that highest posterior density may not be a good criterion for set estimation in scientific inference.

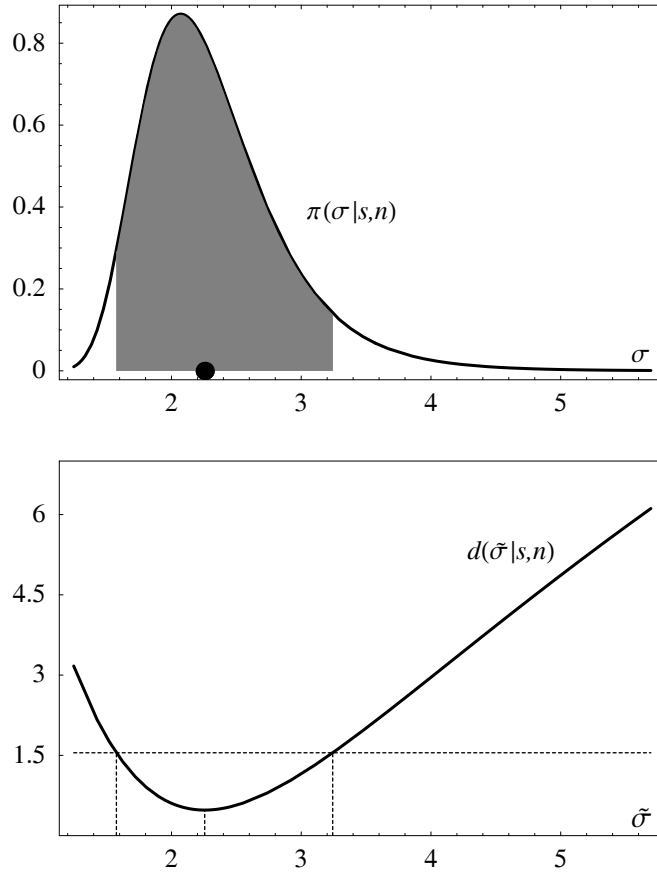
In *one-dimensional* problems, one may define *probability centred* credible intervals, and these are invariant under reparametrization. Indeed, the probability centred  $p$ -credible interval of a real-valued quantity of interest  $\theta$  is defined by the  $(1 - p)/2$  and  $(1 + p)/2$  quantiles of its posterior distribution  $\pi(\theta|\mathbf{x})$ , and this is invariant under one-to-one reparametrizations, since all quantiles are invariant. However, the notion cannot be uniquely extended to multidimensional problems and, even in one-dimensional

problems, their use may be less than satisfactory as, for instance, in those situations where the posterior density is monotonously decreasing within its support.

Whenever a loss structure has been established, foundations dictate that values with smaller expected loss are to be preferred. Thus, for any loss function  $\ell\{\tilde{\theta}, (\theta, \omega)\}$  it is natural to define  $p$ -credible *lowest posterior loss* (LDL) region estimators (Bernardo, 2005b) as those  $p$ -credible regions which contain  $\tilde{\theta}$  values whose expected loss  $l(\tilde{\theta}|\mathbf{x})$  (Eq. 15), is smaller than that of any  $\tilde{\theta}$  values outside the region.

In particular, if the loss function is quadratic, so that

$$\ell\{\tilde{\theta}, (\theta, \lambda)\} = (\tilde{\theta} - \theta)^t(\tilde{\theta} - \theta),$$



**Figure 2:** Reference posterior density of the standard deviation  $\sigma$  of a normal distribution (upper panel), and intrinsic expected loss from using  $\tilde{\sigma}$  as a proxy for  $\sigma$  (lower panel), given a random sample  $\mathbf{x}$  of size  $n = 12$  with standard deviation  $s = 2.071$ . The intrinsic estimate (solid dot) is  $\tilde{\sigma}_{int} = 2.256$ ; the 0.90 intrinsic credible region (shaded region) is  $C_{0.90}^{int} = (1.575, 3.243)$ .

the expected loss is

$$\begin{aligned} l(\tilde{\boldsymbol{\theta}}|\mathbf{x}) &= \int_{\Theta} (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta})^t (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}) \pi(\boldsymbol{\theta}|\mathbf{x}) d\boldsymbol{\theta} \\ &= (\tilde{\boldsymbol{\theta}} - \mathbb{E}[\boldsymbol{\theta}|\mathbf{x}])^t (\tilde{\boldsymbol{\theta}} - \mathbb{E}[\boldsymbol{\theta}|\mathbf{x}]) + \text{Var}[\boldsymbol{\theta}|\mathbf{x}]. \end{aligned}$$

Hence, with quadratic loss, the lowest posterior loss  $p$ -credible region consists of those  $\tilde{\boldsymbol{\theta}}$  values with the smallest Euclidean distance to the posterior mean  $\mathbb{E}[\boldsymbol{\theta}|\mathbf{x}]$ . Notice that these LDL  $p$ -credible regions are *not* invariant under reparametrization.

To obtain LDL invariant region estimators the loss function used must be invariant under one-to-one reparametrization. The arguments mentioned above suggest the use of the intrinsic discrepancy loss. The  $p$ -credible *intrinsic region estimator* is the lowest posterior loss  $p$ -credible region which corresponds to the intrinsic discrepancy loss.

**Definition 5 (Intrinsic region estimator)** Consider once more the parametric family of probability distributions  $\mathcal{F}$  defined by (15). An intrinsic  $p$ -credible region for  $\boldsymbol{\theta}$  is a subset  $C_p^{\text{int}} = C_p^{\text{int}}(\mathbf{x}, \Theta)$  of  $\Theta$  such that

$$\int_{C_p^{\text{int}}} \pi(\boldsymbol{\theta}|\mathbf{x}) d\boldsymbol{\theta} = p, \quad \forall \tilde{\boldsymbol{\theta}}_i \in C_p^{\text{int}}, \forall \tilde{\boldsymbol{\theta}}_j \notin C_p^{\text{int}}, d\{\tilde{\boldsymbol{\theta}}_i|\mathbf{x}\} < d\{\tilde{\boldsymbol{\theta}}_j|\mathbf{x}\},$$

where, again,  $d\{\tilde{\boldsymbol{\theta}}|\mathbf{x}\}$  is the intrinsic expected loss (16).

Intrinsic credible regions are typically unique and, since they are based in the invariant intrinsic discrepancy loss, they are consistent under one-to-one reparametrization. Thus, if  $\boldsymbol{\psi}(\boldsymbol{\theta})$  is a one-to-one function of  $\boldsymbol{\theta}$ , the image  $C_p^{\text{int}}(\mathbf{x}, \Psi) = \boldsymbol{\psi}\{C_p^{\text{int}}(\mathbf{x}, \Theta)\}$  of an intrinsic  $p$ -credible region for  $\boldsymbol{\theta}$  is an intrinsic  $p$ -credible region for  $\boldsymbol{\phi}$ . For more details on intrinsic region estimation, see Bernardo (2005b).

**Example 4 (Normal variance, continued)** Numerical minimization of the intrinsic expected loss (19) in Example 3 immediately yields the intrinsic estimator of the standard deviation  $\sigma$ . This is

$$\sigma^*(\mathbf{x}) = \sigma^*(n, s) = \arg \min_{\tilde{\sigma} > 0} d(\tilde{\sigma}|\mathbf{x}) = 2.256, \quad (20)$$

and it is marked with a solid dot in the top panel of Figure 2. Since intrinsic estimation is invariant, the intrinsic estimates of  $\sigma^2$  or  $\log \sigma$  are respectively  $(\sigma^*)^2$  and  $\log(\sigma^*)$ .

Moreover, the intrinsic  $p$ -credible interval for  $\sigma$  is given by  $C_p^{\text{int}} = (\sigma_0, \sigma_1)$ , where  $\{\sigma_0, \sigma_1\}$  is the unique solution to the equations system

$$\begin{cases} d(\sigma_0|\mathbf{x}) = d(\sigma_1|\mathbf{x}) \\ \int_{\sigma_0}^{\sigma_1} \pi(\sigma|\mathbf{x}) d\sigma = p \end{cases}$$

For instance, with  $p = 0.90$  this yields  $C_{0.90}^{int} = (1.575, 3.243)$ , the shaded region in the top panel of Figure 2. Since intrinsic region estimation is also invariant under reparametrization, the intrinsic  $p$ -credible intervals for  $\sigma^2$  or  $\log \sigma$  will respectively be  $(C_p^{int})^2$  and  $\log(C_p^{int})$ .

### 3 Intrinsic estimation in location-scale models

This section analyses intrinsic point and region estimation of the parameters  $\mu$  and  $\sigma$  (or arbitrary one-to-one functions of these) of variation-independent location-scale models.

#### 3.1 Reference analysis of location-scale models

The likelihood function  $p(\mathbf{x}|\mu, \sigma, f)$  which corresponds to a random sample  $\mathbf{x} = \{x_1, \dots, x_n\}$  from a location-scale model  $m(x|\mu, \sigma, f)$  of the form (2), is given by Equation (3). This will typically have a unique maximum, denoted  $(\hat{\mu}, \hat{\sigma})$ , which not always has a simple analytical expression.

Under appropriate regularity conditions (see, e.g., Bernardo and Smith, 1994, Sec. 5.3 and references therein) the joint posterior distribution of  $\mu$  and  $\sigma$  will be asymptotically normal with mean  $(\hat{\mu}, \hat{\sigma})$  and covariance matrix

$$V(\hat{\mu}, \hat{\sigma}, n) = n^{-1} F^{-1}(\hat{\mu}, \hat{\sigma}) = (\hat{\sigma}^2/n) A^{-1}(f) \quad (21)$$

where  $F(\mu, \sigma)$  is Fisher's information matrix which, in location-scale models, is always of the form  $F(\mu, \sigma) = \sigma^{-2}A(f)$ , where  $A(f)$  is a  $2 \times 2$  matrix which depends on the probability density  $f(\cdot)$ , but not on  $\mu$  or  $\sigma$ . As a consequence, if the parameter of interest is either  $\mu$  (or a one-to-one function of  $\mu$ ) or  $\sigma$  (or a one-to-one function of  $\sigma$ ) with independent variation, then (Fernández and Steel, 1999, Th. 1), under regularity conditions sufficient to guarantee posterior asymptotic normality, and variation independence of  $\mu$  and  $\sigma$ , the joint reference prior is independent of the function  $f(\cdot)$  and, in terms of  $\mu$  and  $\sigma$ , is given by

$$\pi(\mu) \pi(\sigma | \mu) = \pi(\sigma) \pi(\mu | \sigma) = \sigma^{-1}. \quad (22)$$

Using Bayes theorem, the corresponding joint reference posterior is

$$\pi(\mu, \sigma | \mathbf{x}) = \sigma^{-(n+1)} \prod_{j=1}^n f\{\sigma^{-1}(x_j - \mu)\}, \quad (23)$$

which is typically proper for all  $n \geq 2$ . In particular, it is proper for all  $n \geq 2$  whenever

$f(\cdot)$  is either a standard normal or a scale mixture of standard normals, what includes Student models.

In the normal case, with  $f(x) = N(x|0, 1)$ , the joint posterior (23) becomes

$$\pi(\mu, \sigma | \bar{x}, s, n) = N(\mu | \bar{x}, \sigma / \sqrt{n}) \pi(\sigma | s, n),$$

where  $\pi(\sigma | s, n)$  is the square root inverted gamma given by (5). The corresponding marginal reference posterior of the precision  $\lambda = \sigma^{-2}$  is found to be  $\pi(\lambda | \mathbf{x}) = \text{Ga}(\lambda | (n-1)/2, (ns^2)/2)$  and, thus,

$$E[\lambda | \mathbf{x}] = \frac{n-1}{n s^2}, \quad \text{Var}[\lambda | \mathbf{x}] = \frac{2(n-1)}{n^2 s^4}. \quad (24)$$

The marginal reference posterior of  $\mu$  is the Student distribution

$$\pi(\mu | \mathbf{x}) = \text{St}\left(\mu | \bar{x}, \frac{s}{\sqrt{n-1}}, n\right) \propto \left(1 + \frac{(\mu - \bar{x})^2}{s^2}\right)^{-n/2}. \quad (25)$$

For details see, for instance, Bernardo and Smith (1994, Sec. 5.4).

### 3.2 Intrinsic estimation of the normal mean

As stated in Example 2 (Eq. 11), the directed divergence  $\kappa\{\mu_j, \sigma_j | \mu_i, \sigma_i\}$ , of  $N(x | \mu_j, \sigma_j)$  from  $N(x | \mu_i, \sigma_i)$ , is

$$\kappa\{\mu_j, \sigma_j | \mu_i, \sigma_i\} = \frac{1}{2} \left\{ \frac{\sigma_i^2}{\sigma_j^2} - 1 - \log \frac{\sigma_i^2}{\sigma_j^2} + \frac{(\mu_i - \mu_j)^2}{\sigma_j^2} \right\}.$$

As a function of  $\tilde{\sigma}$ , the directed divergence  $\kappa\{\tilde{\mu}, \tilde{\sigma} | \mu, \sigma\}$  is minimized when  $\tilde{\sigma}^2$  takes the value  $\tilde{\sigma}_{min}^2 = (\mu - \mu_i)^2 + \sigma^2$ , and substitution yields

$$\kappa\{\tilde{\mu}, \tilde{\sigma}_{min} | \mu, \sigma\} = \frac{1}{2} \log \left[ 1 + \frac{(\mu - \tilde{\mu})^2}{\sigma^2} \right].$$

Similarly, the directed divergence  $\kappa\{\mu, \sigma | \tilde{\mu}, \tilde{\sigma}\}$  is minimized, as a function of  $\tilde{\sigma}$ , when  $\tilde{\sigma} = \sigma$ , and substitution now yields

$$\kappa\{\mu, \sigma | \tilde{\mu}, \sigma\} = \frac{1}{2} \frac{(\mu - \tilde{\mu})^2}{\sigma^2}.$$

Hence, making use of Theorem 1 and the fact that, for all  $x > 0$ ,  $\log(1+x) \leq x$ , the intrinsic discrepancy loss  $\delta\{\tilde{\mu}, (\mu, \sigma)\}$  from using  $\tilde{\mu}$  as a proxy for  $\mu$  with a normal sample of size  $n$  is

$$\delta\{\tilde{\mu}, (\mu, \sigma)\} = \delta\{\theta^2\} = \frac{n}{2} \log\left[1 + \frac{\theta^2}{n}\right], \quad \theta = \theta(\tilde{\mu}, \mu, \sigma) = \frac{\mu - \tilde{\mu}}{\sigma/\sqrt{n}}, \quad (26)$$

which only depends on the number  $\theta$  of standard deviations which separate  $\tilde{\mu}$  from  $\mu$ . Figure 3 represents the intrinsic loss function (26), as a function of  $\theta$ , for several values of  $n$ .

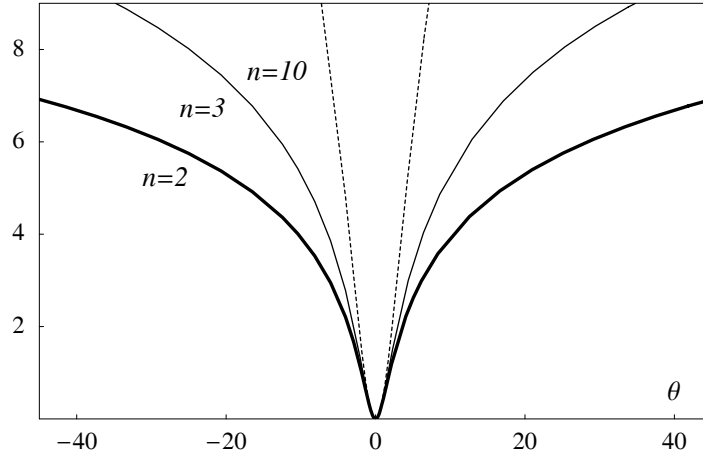
As one might expect,  $\delta\{\tilde{\mu}, (\mu, \sigma)\}$  increases with  $|\theta|$ . The dependence is essentially quadratic in a neighbourhood of zero, but shows a very reasonable concavity in regions where  $|\theta|$  is large.

Using Definition 3, the intrinsic discrepancy reference expected loss  $d(\tilde{\mu} | \mathbf{x})$  may be written in terms of the reference posterior of  $\theta$ ; indeed,

$$\begin{aligned} d(\tilde{\mu} | \mathbf{x}) &= \int_0^\infty \int_{-\infty}^\infty \delta\{\tilde{\mu}, (\mu, \sigma)\} \pi(\mu, \sigma | \mathbf{x}) d\mu d\sigma \\ &= \int_0^\infty \frac{n}{2} \log\left[1 + \frac{\theta^2}{n}\right] \pi(\theta | \mathbf{x}) d\theta. \end{aligned} \quad (27)$$

But  $\theta = (\tilde{\mu} - \mu)/(\sigma/\sqrt{n})$  may be written as  $a + \beta$  where, as a function of  $\mu$  and  $\sigma$ ,  $\beta = (\mu - \bar{x})/(\sigma/\sqrt{n})$  has a standard normal reference posterior, and  $a$  is the constant  $a = (\bar{x} - \tilde{\mu})/(\sigma/\sqrt{n})$ . Hence, the conditional posterior distribution of  $\theta^2$  given  $\sigma$  is noncentral  $\chi^2$  with one degree of freedom and non centrality parameter  $a^2$ ,

$$\pi(\theta^2 | \mathbf{x}, \sigma) = \chi^2(\theta^2 | 1, a^2), \quad a^2 = n \frac{(\bar{x} - \tilde{\mu})^2}{\sigma^2}.$$



**Figure 3:** Intrinsic discrepancy loss for estimation of the normal mean as a function of the number  $\theta = (\tilde{\mu} - \mu)/(\sigma/\sqrt{n})$  of standard deviations which separates  $\tilde{\mu}$  from  $\mu$ , for  $n = 2$ ,  $n = 3$ , and  $n = 10$ .



It follows that the intrinsic expected loss  $d(\tilde{\mu} | \mathbf{x})$  only depends on  $\tilde{\mu}$  through  $(\bar{x} - \tilde{\mu})^2$ , and increases with  $(\bar{x} - \tilde{\mu})^2$ ; therefore, the intrinsic estimator of  $\mu$  is

$$\tilde{\mu}_{int}(\mathbf{x}) = \arg \min_{\tilde{\mu} \in \mathcal{R}} d(\tilde{\mu} | \mathbf{x}) = \arg \min_{\tilde{\mu} \in \mathcal{R}} (\bar{x} - \tilde{\mu})^2 = \bar{x}.$$

Moreover,  $d(\tilde{\mu} | \mathbf{x})$  is symmetric around  $\bar{x}$  and, hence, all intrinsic credible regions must be centered at  $\bar{x}$ . In view of (25), this implies that the intrinsic the  $p$ -credible regions are just the usual Student- $t$  HPD  $p$ -credible intervals

$$C_p^{int}(\mathbf{x}) = \left\{ \tilde{\mu}; \bar{x} - q_{p,n} s / \sqrt{n-1} \leq \tilde{\mu} \leq \bar{x} + q_{p,n} s / \sqrt{n-1} \right\}, \quad (28)$$

where  $q_{p,n}$  is the  $(p+1)/2$  quantile of a standard Student- $t$  with  $n-1$  degrees of freedom.

It immediately follows from (28) that  $C_p^{int}$  consist of the set of  $\tilde{\mu}$  values such that  $(\bar{x} - \tilde{\mu})/(s/\sqrt{n-1})$  belongs to a probability  $p$  centred interval of a standard Student- $t$  with  $n-1$  degrees of freedom. But, as a function of the data  $\mathbf{x}$ , the sampling distribution of

$$t(\mathbf{x}) = (\bar{x} - \mu)/(s/\sqrt{n-1}) \quad (29)$$

is also a standard Student- $t$  with  $n-1$  degrees of freedom. Hence, for all sample sizes, the *expected coverage under sampling* of the  $p$ -credible intervals (28) is *exactly*  $p$ , and the intrinsic credible regions are exact frequentist confidence intervals.

A simple asymptotic approximation to  $d(\tilde{\mu} | \mathbf{x})$ , which provides a direct measure in a log-likelihood ratio scale of the expected loss associated to the use of  $\tilde{\mu}$ , may easily be obtained. Indeed, a variation of the delta method shows that, under appropriate regularity conditions, the expectation of some function  $y = g(x)$  of a random quantity  $x$  with mean  $\mu_x$  and variance  $\sigma_x^2$  may be approximated by

$$E[g(x)] \approx g \left[ \mu_x + \frac{\sigma_x^2}{2} \frac{g''(\mu_x)}{g'(\mu_x)} \right]. \quad (30)$$

On the other hand, the conditional posterior mean of  $\theta^2$  is  $1 + a^2$ , and its conditional posterior variance is  $2 + 4a^2$ ; but  $E[\sigma^{-2} | \mathbf{x}] = E[\lambda | \mathbf{x}] = (n-1)/(ns^2)$  (Eq. 24) and hence, the unconditional posterior mean and variance of  $\theta^2(\tilde{\mu})$  are, respectively,

$$E[\theta^2 | \mathbf{x}] = 1 + t^2, \quad \text{Var}[\theta^2 | \mathbf{x}] = 2 + 4t^2,$$

both functions of the conventional  $t$  statistic (29). Using these in (30) to approximate the posterior expectation of  $\log(1 + \theta^2/n)$  required in (27) yields

$$d(\tilde{\mu} | \mathbf{x}) \approx \frac{n}{2} \log \left[ 1 + \frac{1}{n} \frac{n(1+t^2) + t^4}{n+t^2+1} \right]. \quad (31)$$

Progressively cruder, but simpler approximations are

$$d(\tilde{\mu} | \mathbf{x}) \approx \frac{n}{2} \log \left[ 1 + \frac{1}{n} (1 + t^2) \right] \approx \frac{1}{2} (1 + t^2). \quad (32)$$

Thus, for large  $n$ , the intrinsic expected loss  $d(\tilde{\mu} | \mathbf{x})$  is essentially quadratic in the number  $t = (\bar{x} - \tilde{\mu})/(s/\sqrt{n-1})$  of standard deviations which separate  $\bar{x}$  from  $\tilde{\mu}$ . Summarizing, we have thus established

**Theorem 2 (Intrinsic estimation of the Normal mean)** *Let  $\mathbf{x}$  be a random sample of size  $n$  from  $N(x|\mu, \sigma)$ , with mean and variance  $\bar{x}$  and  $s^2$ , and let  $t = \sqrt{n-1}(\bar{x} - \tilde{\mu})/s$  be the conventional  $t$  statistic.*

(i) *The intrinsic point estimator of  $\mu$  is  $\tilde{\mu}_{int}(\mathbf{x}) = \bar{x}$ .*

(ii) *The unique  $p$ -credible intrinsic region for  $\mu$  is the probability centred interval*

$$C_p^{int}(\mathbf{x}) = \bar{x} \pm q_{p,n} s / \sqrt{n-1},$$

*where  $q_{p,n}$  is the  $(p+1)/2$  quantile of a standard Student- $t$  distribution with  $n-1$  degrees of freedom. For all sample sizes, the frequentist coverage of  $C_p^{int}(\mathbf{x})$  is exactly  $p$ .*

(iii) *The expected intrinsic loss associated to the use of  $\tilde{\mu}$  as a proxy for  $\mu$  is*

$$d(\tilde{\mu} | \mathbf{x}) \approx \frac{n}{2} \log \left[ 1 + \frac{1}{n} \frac{n(t^2 + 1) + t^4}{n + t^2 + 1} \right] \approx \frac{n}{2} \log \left[ 1 + \frac{1}{n} (1 + t^2) \right].$$

As a numerical illustration, a random sample of size  $n = 25$  was generated from a standard normal, yielding  $\bar{x} = -0.162$  and  $s = 0.840$ . The intrinsic estimator is  $\mu^* = \bar{x} = -0.162$  and the 0.99-intrinsic credible region is the interval  $[-0.642, 0.318]$ . The exact value of the expected intrinsic loss  $d(1/3 | \mathbf{x})$ , computed from (27) by numerical integration, is 3.768, while (31) and the two approximations in (32) respectively yield 3.781, 3.970 and 4.673. Hence, the observed data may be expected to be about  $\exp(3.768) \approx 43$  times more likely under the true value of  $\mu$  than under the closest normal model with  $\mu = 1/3$ , suggesting that the value  $\mu = 1/3$  is hardly compatible with the observed data.

### 3.3 Intrinsic estimation of the normal variance

It has already been established (Example 2, Eq. 13) that the intrinsic discrepancy loss from using  $\tilde{\sigma}^2$  as a proxy for  $\sigma^2$  is

$$\delta_x\{\tilde{\sigma}^2, (\mu, \sigma)\} = \delta_x\{\phi\} = \begin{cases} \frac{n}{2} g(\phi) & \text{if } \phi < 1, \\ \frac{n}{2} g(1/\phi) & \text{if } \phi \geq 1, \end{cases} \quad (33)$$

where  $g(t) = (t - 1) - \log t$ , and  $\phi = \tilde{\sigma}^2/\sigma^2$ , and that this is also the intrinsic loss  $\delta_x\{\tilde{\psi}, (\mu, \sigma)\}$  from using  $\tilde{\psi}$  as a proxy for  $\psi$  for any one-to-one function  $\psi(\sigma^2)$  of  $\sigma^2$ . Moreover, the reference posterior of  $\phi$  is the gamma distribution of Eq. 18. Hence, the intrinsic estimator of  $\sigma^2$  is

$$\tilde{\sigma}_{int}^2(\mathbf{x}) = \arg \min_{\tilde{\sigma}^2 > 0} \int_0^\infty \delta_x(\phi) \text{Ga}\left(\phi \mid \frac{n-1}{2}, \frac{n s^2}{2\tilde{\sigma}^2}\right) d\phi,$$

where  $\delta_x(\phi)$  is given by (33). Moreover, it immediately follows from (18) that, as a function of  $\sigma$ , the reference posterior distribution of  $\tau = n s^2/\sigma^2$  is

$$\pi(\tau | \mathbf{x}) = \pi(\tau | n) = \chi^2(\tau | n - 1) \quad (34)$$

a central  $\chi^2$  with  $n - 1$  degrees of freedom; but  $\phi = c \tau/n$ , with  $c = \tilde{\sigma}^2/s^2$  and, therefore, the expected posterior loss from using  $\tilde{\sigma}$  may further be written as

$$d(\tilde{\sigma}^2 | s^2, n) = d(c | n) = \int_0^\infty \delta\left(\frac{c \tau}{n}\right) \chi^2(\tau | n - 1) d\tau, \quad c = \tilde{\sigma}^2/s^2. \quad (35)$$

Thus, the intrinsic estimator of the normal variance is an affine equivariant estimator of the form

$$\sigma_{int}^2(s, n) = c_n^* s^2, \quad c_n^* > 0, \quad (36)$$

where  $c_n^*$  is the value of  $c$  which minimizes  $d(c | n)$  in (35). The exact value of  $c_n^*$  may be numerically found by one-dimensional numerical integration, followed by numerical optimization. The first row of Table 1 displays the exact values of  $c_n^*$  for several sample sizes. However, good analytical approximations for  $c_n^*$  may be obtained.

We first consider a general approximation method. Let  $\omega$  be a particular parametrization of the problem, and consider a (variance stabilizing) *reference reparametrization*

**Table 1:** Exact and alternative approximate values for the intrinsic point estimator of the normal variance  $\sigma_{int}^2 = c_n^* s^2$ .

$n$	2	3	4	5	10	20	50
$c_n^*$	4.982	2.347	1.803	1.569	1.231	1.106	1.041
$\left(\frac{n}{n-1}\right)^2$	4.000	2.250	1.778	1.563	1.235	1.108	1.041
$\frac{n}{n-1} e^{1/(n-1)}$	5.437	2.473	1.861	1.605	1.242	1.110	1.041
$\frac{n}{n-2}$	—	3.000	2.000	1.667	1.250	1.111	1.042

$\phi(\omega)$  defined as one with a uniform reference prior. This is given by any solution to the differential equation  $\phi'(\omega) = \pi(\omega)$ , where  $\pi(\omega)$  is the marginal reference prior for  $\omega$ . Under regularity conditions, the sampling distribution of  $\phi(\hat{\omega})$ , where  $\hat{\omega} = \hat{\omega}(\mathbf{x})$  is the MLE of  $\omega$ , and the reference posterior of  $\phi(\omega)$ , are both asymptotically normal. Using these approximations, the intrinsic expected loss from using  $\tilde{\omega}$  is found to be (Bernardo, 2005b, Theo. 4.1)

$$d(\tilde{\omega} | \mathbf{x}) \approx \frac{n}{2} \{ \sigma_\phi^2 + [\mu_\phi - \phi(\tilde{\omega})]^2 \}. \quad (37)$$

where  $\mu_\phi$  and  $\sigma_\phi^2$  are respectively the posterior mean and posterior variance of  $\phi = \phi(\omega)$ . This is minimized by  $\phi(\tilde{\omega}) = \mu_\phi = E[\phi | \mathbf{x}]$ . Hence, in terms of any reference parametrization  $\phi$ , the intrinsic point estimate is approximately the posterior mean  $\mu_\phi$  and, by invariance, the intrinsic estimator of any one-to-one function,  $\psi = \psi(\phi)$  is approximately given by  $\tilde{\psi}_{int} = \psi(\mu_\phi)$ . Thus,

$$\tilde{\phi}_{int}(\mathbf{x}) \approx \mu_\phi = E[\phi | \mathbf{x}], \quad \tilde{\omega}_{int}(\mathbf{x}) \approx \phi^{-1}\{\mu_\phi\}. \quad (38)$$

Under regularity conditions (see, e.g., Schervish, 1995, Sec. 7.1.3) the delta method may be used to obtain simple approximations to the posterior moments of  $\phi$  in terms of those of  $\omega$ , namely

$$\mu_\phi \approx \phi\{\mu_\omega\} + \sigma_\omega^2 \phi''\{\mu_\omega\}/2, \quad (39)$$

$$\sigma_\phi^2 \approx \sigma_\omega^2 [\phi'\{\mu_\omega\}]^2. \quad (40)$$

Substitution into (38) and (37) respectively provide useful approximations to the intrinsic point estimator of  $\omega$ , and to the expected loss from using  $\tilde{\omega}$  as a proxy for  $\omega$ .

In the particular case of the normal variance, it is convenient to start from the parametrization in terms of the precision  $\lambda = \sigma^{-2}$ , whose posterior moments have simple expressions. Since reference priors are consistent under reparametrization, the reference prior for  $\lambda$  is  $\pi(\lambda) = \pi(\sigma)|\partial\sigma/\partial\lambda| \propto \lambda^{-1}$  and, therefore, a reference parametrization is

$$\phi = \phi(\lambda) = \int^\lambda \pi(\lambda) d\lambda = \int^\lambda \lambda^{-1} d\lambda = \log \lambda.$$

Notice that the reference prior of  $\phi = \log \lambda$  is indeed uniform, as it is the case for the logarithm of any other power of  $\sigma$ . Using (39) and (40) with the first posterior moments of  $\lambda$ , given in (24), yields

$$\tilde{\phi}_{int}(\mathbf{x}) \approx \mu_\phi = E[\log \lambda | \mathbf{x}] \approx \log\left(\frac{n-1}{n s^2}\right) - \frac{1}{n-1}, \quad (41)$$

$$\sigma_\phi^2 = \text{Var}[\log \lambda | \mathbf{x}] \approx \frac{2}{n-1}. \quad (42)$$

By invariance, (41) directly provides an approximation to the intrinsic estimator of the variance. This has the form of a modified version of the conventional unbiased estimator  $\tilde{\sigma}_{n-1}^2$ ; indeed, since  $\sigma^2 = e^{-\phi}$ ,

$$\tilde{\sigma}_{int}^2(\mathbf{x}) = e^{-\tilde{\phi}_{int}(\mathbf{x})} \approx \frac{n s^2}{n-1} e^{\frac{1}{n-1}} = \tilde{\sigma}_{n-1}^2 e^{\frac{1}{n-1}},$$

which, as shown in the third row of Table 1, provides good approximations, even for small values of  $n$ .

A better analytical approximation to the intrinsic estimator of the normal variance may be obtained making use of the particular features of this example. This is done by separately minimizing the expected value of each of the two functions which enter the definition of the intrinsic discrepancy loss  $\delta_x\{\phi\}$ , and using the arithmetic mean of the corresponding results.

Indeed, the delta method may be used to approximate both  $E[g(c\tau/n)]$  and  $E[g(n/(c\tau))]$  in terms  $E[\tau|n] = n-1$  and  $\text{Var}[\tau|n] = 2(n-1)$ . The approximation to  $E[g(c\tau/n)]$  is minimized by  $\hat{c}_{1n}^* = n/(n-1)$ , while the approximation to  $E[g(n/(c\tau))]$  is minimized by  $\hat{c}_{2n}^* = n(n+1)/(n-1)^2$ . As one would expect, their average,

$$\hat{c}_n^* = \frac{\hat{c}_{1n}^* + \hat{c}_{2n}^*}{2} = \left(\frac{n}{n-1}\right)^2 = \frac{n}{n-(2-n^{-1})}, \quad (43)$$

provides a good approximation to the value  $c_n^*$  which minimizes (35). As shown in the second row of Table 1, the approximation remains good even for small values of  $n$ . Combination of (36) and (43) establishes that, for all but very small  $n$  values,

$$\tilde{\sigma}_{int}^2(\mathbf{x}) = \tilde{\sigma}_{int}^2(s, n) \approx \left(\frac{n}{n-1}\right)^2 s^2. \quad (44)$$

In view of the second expression for  $\hat{c}_n^*$  in (43), a cruder approximation is given by  $\tilde{\sigma}_{int}^2 \approx \tilde{\sigma}_{n-2}^2$ . This is larger than the MLE  $\hat{\sigma}^2 = s^2$  (which divides by  $n$  the sum of squares), and also larger than the conventional unbiased estimate of the variance  $\tilde{\sigma}_{n-1}^2$  (which divides by  $n-1$ ). Notice that numerical differences between intrinsic and conventional estimators may be large for small values of  $n$ . In particular, with only two observations  $\{x_1, x_2\}$ , the intrinsic estimator of the variance is  $\tilde{\sigma}_{int}^2(2, s^2) \approx 5 s^2 = 5(x_1 - x_2)^2/4$ ; this is 2.5 times larger than the unbiased estimator,  $(x_1 - x_2)^2/2$  in this case, which (with good reason) is generally considered to be too small.

As shown by (35), the expected intrinsic loss  $d(\tilde{\sigma}^2 | n, s^2)$  of any affine equivariant estimator of the variance  $\tilde{\sigma}^2 = k_n s^2$ , is actually independent of  $s^2$  and only depends on the sample size  $n$ . Moreover, it is easily verified that the expected intrinsic loss  $d(\tilde{\sigma}^2 | n, s^2)$  is precisely equal to the *frequentist risk* associated to the intrinsic discrepancy loss,

$$r(\tilde{\sigma}_i^2 | n, \sigma^2) = \int_0^\infty \delta\{\tilde{\sigma}_i^2(s^2), \sigma^2\} p(s^2 | n, \sigma^2) ds^2.$$

Thus, under intrinsic discrepancy loss, the intrinsic estimator  $\tilde{\sigma}_{int}^2$  dominates all affine equivariant estimators. For details, see Bernardo (2006).

Region estimation is now considered. As described in Example 4, the intrinsic  $p$ -credible region for  $\sigma$  is the unique solution  $C_p^{int} = \{\sigma_0, \sigma_1\}$  to the equations system

$$\left\{ d(\sigma_0^2 | \mathbf{x}) = d(\sigma_1^2 | \mathbf{x}), \quad \int_{\sigma_0}^{\sigma_1} \pi(\sigma | \mathbf{x}) d\sigma = p \right\}.$$

Using (34), this may equivalently be written in terms of  $\tau = n s^2 / \sigma^2$  as

$$\left\{ d(\sigma_0^2 | \mathbf{x}) = d(\sigma_1^2 | \mathbf{x}), \quad \int_{ns^2/\sigma_1}^{ns^2/\sigma_0} \chi(\tau | n-1) d\tau = p \right\}. \quad (45)$$

Thus, the unique  $p$ -credible intrinsic region for  $\sigma^2$  is the interval

$$C_p^{int}(\mathbf{x}) = \left\{ \frac{n s^2}{\chi_{n-1}^2(1-\alpha)}, \frac{n s^2}{\chi_{n-1}^2(1-p-\alpha)} \right\} \quad (46)$$

where  $\chi_{n-1}^2(q)$  is the  $q$  quantile of a  $\chi_{n-1}^2$  distribution, and  $\alpha$  is the solution to the equation

$$d\left(\frac{n s^2}{\chi_{n-1}^2(1-\alpha)} | \mathbf{x}\right) = d\left(\frac{n s^2}{\chi_{n-1}^2(1-p-\alpha)} | \mathbf{x}\right). \quad (47)$$

By invariance, this provides the intrinsic  $p$ -credible region of any one-to-one function of  $\sigma^2$ .

As a function of the data  $\mathbf{x}$ , the sampling distribution of  $n s^2 / \sigma^2$  is also a  $\chi^2$  with  $n-1$  degrees of freedom. Hence, for all sample sizes, the expected coverage under sampling of the  $p$ -credible intervals (46) is *exactly*  $p$ .

Using (35) to evaluate expected losses, the exact solution to equation (47) may easily be obtained by numerical methods. However, good analytical approximations may be obtained.

Working again in terms of the reference parametrization for this problem,  $\phi = \log \lambda = -2 \log \sigma$ , and using (37), (39) and (40), the expected loss from using  $\tilde{\phi}$  as a proxy for  $\phi$  is approximately

$$d(\tilde{\phi} | \mathbf{x}) \approx \frac{n}{2} \left[ \frac{2}{n-1} + (\tilde{\phi}_{int} - \tilde{\phi})^2 \right]. \quad (48)$$

But this is symmetric around  $\tilde{\phi}_{int} = \log(\tilde{\lambda}_{int}) = -\log(\sigma_{int}^2)$  and therefore, to keep those  $\tilde{\phi}$  points with smaller expected loss, any intrinsic credible region for  $\phi = \log \lambda$  must be (approximately) centered at  $\tilde{\phi}_{int}$ . Thus, using (42) and (44) this will be of the form

$$C_p^{int}(\mathbf{x}, \Phi) \approx \tilde{\phi}_{int} \pm \alpha_{pn} \sigma_\phi \approx \log \left[ \left( \frac{n-1}{n} \right)^2 \frac{1}{s^2} \right] \pm \gamma_{pn} \sqrt{\frac{2}{n-1}} \quad (49)$$

where  $\gamma_{pn}$  is the solution to the equation

$$\int_{\tilde{\phi}_{int} - \gamma_{pn} \sigma_\phi}^{\tilde{\phi}_{int} + \gamma_{pn} \sigma_\phi} \pi(\phi | \mathbf{x}) d\phi = p,$$

or, equivalently since

$$\begin{aligned} \tau = n s^2 \lambda &= n s^2 e^\phi, \\ n s^2 e^{\tilde{\phi}_{int} \pm \gamma_{pn} \sigma_\phi} &= \frac{(n-1)^2}{n} \exp \left[ \pm \gamma_{pn} \sqrt{\frac{2}{n-1}} \right], \\ \pi(\tau | \mathbf{x}) &= \chi^2(\tau | n-1), \end{aligned}$$

$\gamma_{pn}$  is the unique solution to the equation

$$F_{n-1} \left\{ \frac{(n-1)^2}{n} e^{+\gamma_{pn} \sqrt{\frac{2}{n-1}}} \right\} - F_{n-1} \left\{ \frac{(n-1)^2}{n} e^{-\gamma_{pn} \sqrt{\frac{2}{n-1}}} \right\} = p, \quad (50)$$

where  $F_\nu$  is the cumulative distribution function of a  $\chi_\nu^2$  distribution.

A numerical solution to (50) is immediately found with standard statistical software. However, a simple analytical approximation may be derived using the fact that the reference posterior distribution of  $\phi = \log \lambda$  becomes approximately normal (at a faster rate than any other simple function of  $\sigma$ ) as the sample size  $n$  increases. Using this approximation and (49), the  $p$ -credible intrinsic region for  $\phi$  is approximated by the interval

$$C_p^{int}(\mathbf{x}, \Phi) \approx \log \left[ \left( \frac{n-1}{n} \right)^2 \frac{1}{s^2} \right] \pm q_p \sqrt{\frac{2}{n-1}} \quad (51)$$

where  $q_p$  is the  $(p+1)/2$  quantile of a standard normal distribution. By invariance, the  $p$ -credible intrinsic region for the variance  $\sigma^2 = e^{-\phi}$  will be approximated by

$$\left\{ s^2 \left( \frac{n}{n-1} \right)^2 e^{-\gamma_{np} \sqrt{\frac{2}{n-1}}}, s^2 \left( \frac{n}{n-1} \right)^2 e^{+\gamma_{np} \sqrt{\frac{2}{n-1}}} \right\} \quad (52)$$

where  $\gamma_{np}$  is the solution to (50) which, as  $n$  increases, converges to  $q_p$ , the  $(p+1)/2$  quantile of an standard normal distribution.

Summarizing, we have thus established

**Theorem 3 (Intrinsic estimation of the normal variance)** Let  $\mathbf{x}$  be a random sample of size  $n$  from  $N(x | \mu, \sigma)$ , with variance  $s^2$ .

(i) The intrinsic point estimator  $\tilde{\sigma}_{int}^2(\mathbf{x})$  of  $\sigma^2$  is

$$\begin{aligned}\tilde{\sigma}_{int}^2(\mathbf{x}) &= \arg \min_{\tilde{\sigma}^2 > 0} d(\tilde{\sigma}^2 | \mathbf{x}), \\ d(\tilde{\sigma}^2 | \mathbf{x}) &= \frac{n}{2} \int_0^\infty \delta\left(\frac{\tilde{\sigma}^2 \tau}{n s^2}\right) \chi^2(\tau | n-1) d\tau, \\ \delta\{\theta\} &= \min\{g(\theta), g(1/\theta)\}, \quad g(\theta) = (\theta - 1) - \log \theta, \\ \tilde{\sigma}_{int}^2(\mathbf{x}) &\approx \left(\frac{n}{n-1}\right)^2 s^2.\end{aligned}$$

The intrinsic point estimator  $\tilde{\sigma}_{int}(\mathbf{x})$  is the Bayes estimator with respect to the intrinsic discrepancy loss. Besides, it has smaller frequentist risk with respect to this loss than any other affine equivariant estimator.

(ii) The unique  $p$ -credible intrinsic region  $C_p^{int}(\mathbf{x})$  for  $\sigma^2$  is the interval

$$C_p^{int}(\mathbf{x}) = \{a(\alpha, \mathbf{x}), b(\alpha, p, \mathbf{x})\} = \left\{ \frac{n s^2}{\chi_{n-1}^2(1-\alpha)}, \frac{n s^2}{\chi_{n-1}^2(1-p-\alpha)} \right\},$$

where  $\chi_{n-1}^2(q)$  is the  $q$  quantile of a  $\chi_{n-1}^2$  distribution, and  $\alpha$  is the solution to the equation  $d\{a(\alpha, \mathbf{x}) | \mathbf{x}\} = d\{b(\alpha, p, \mathbf{x}) | \mathbf{x}\}$ . For all sample sizes, the frequentist coverage of  $C_p^{int}(\mathbf{x})$  is exactly  $p$ . Moreover,

$$C_p^{int}(\mathbf{x}) \approx s^2 \left(\frac{n}{n-1}\right)^2 \left\{ e^{-\gamma_{np}} \sqrt{\frac{2}{n-1}}, e^{+\gamma_{np}} \sqrt{\frac{2}{n-1}} \right\}$$

where  $\gamma_{np}$  is the solution to the equation

$$F_{n-1}\left\{\frac{(n-1)^2}{n} e^{+\gamma_{np}} \sqrt{\frac{2}{n-1}}\right\} - F_{n-1}\left\{\frac{(n-1)^2}{n} e^{-\gamma_{np}} \sqrt{\frac{2}{n-1}}\right\} = p.$$

and  $F_v$  is the cumulative distribution function of a  $\chi_v^2$  distribution. As  $n$  increases,  $\gamma_{np}$  converges to the  $(p+1)/2$  normal quantile.

(iii) The expected intrinsic loss associated to the use of  $\tilde{\sigma}^2$  is

$$d(\tilde{\sigma}^2 | s^2, n) \approx \frac{n}{2} \left[ \frac{2}{n-1} + \left( \log \frac{1}{\sigma_{int}^2} - \log \frac{1}{\tilde{\sigma}^2} \right)^2 \right],$$

with  $\tilde{\sigma}_{int}^2(\mathbf{x}) \approx n^2 s^2 / (n-1)^2$ .

For the numerical illustration considered in Example 4 (where the sample size was only  $n = 12$ ), the approximation (44) to the intrinsic estimate of  $\sigma^2$  yields  $\tilde{\sigma}_{int}^2 \approx 5.104$ . The approximation (52) to the intrinsic 0.90-credible region yields (2.480, 11.507) using the exact solution  $\gamma_{np} = 1.693$  to equation (50), and (2.531, 11.272) using the



corresponding normal approximation  $\gamma_{np} \approx q_{0.90} = 1.645$ . These approximations may be compared with the exact values  $\tilde{\sigma}_{int}^2 = 5.090$  and  $C_{0.90}^{int} = (2.481, 11.717)$  numerically found in Example 4.

### 3.4 Intrinsic estimation of the Cauchy scale parameter

We finally consider an example where no analytical expressions are possible. With the use of increasingly complex statistical models, this is fast becoming the norm, rather than the exception, in statistical practice.

Let  $\mathbf{x} = \{x_1, \dots, x_n\}$  be a random sample from a Cauchy distribution  $\text{Ca}(x|0, \sigma)$ , centered at zero with unknown scale parameter  $\sigma$ , so that the likelihood function is

$$p(\mathbf{x}|\sigma) = \prod_{i=1}^n \text{Ca}(x_i|0, \sigma) \propto \sigma^{-n} \prod_{i=1}^n \left(1 + \frac{x_i^2}{\sigma^2}\right)^{-1}.$$

The Cauchy distribution does not belong to the exponential family and, therefore, there is no sufficient statistic of finite dimension. There is no analytical expression for the MLE  $\hat{\sigma}$  of the unknown parameter. Fisher information function is  $n/(2\sigma^2)$  and, therefore, the posterior distribution of  $\sigma$  will be asymptotically normal,  $N(\sigma|\hat{\sigma}, \sqrt{2} \hat{\sigma}/\sqrt{n})$ .

Since this is a scale model, the reference prior is  $\pi(\sigma) = \sigma^{-1}$  and, using Bayes theorem, the reference posterior is

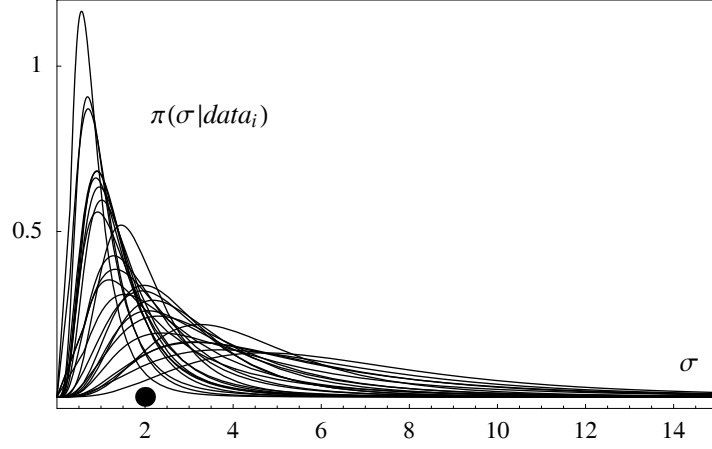
$$\pi(\sigma|\mathbf{x}) = \frac{\sigma^{-1} p(\mathbf{x}|\sigma)}{\int_0^\infty \sigma^{-1} p(\mathbf{x}|\sigma) d\sigma}, \quad (53)$$

which may easily be numerically computed. It may be verified that, provided the data  $\mathbf{x}$  contain at least two different observations,  $\pi(\sigma|\mathbf{x})$  has a gamma-like shape with a unique mode.

Figure 4 represents the reference posteriors of  $\sigma$  which correspond to a set of 25 random samples of size  $n = 12$ , which were all generated from a Cauchy distribution  $\text{Ca}(x|0, 2)$ . This may be seen as a graphical representation of the *sampling distribution* of  $\pi_\sigma(\cdot|\mathbf{x})$ , the reference posterior of  $\sigma$ , given  $\sigma = 2$  and  $n = 12$ . Notice that, although all these posteriors contain indeed the true value  $\sigma = 2$  from which the samples have been simulated (marked in the figure with a solid dot), the variability is very large.

The logarithmic divergence  $\kappa\{\sigma_2|\sigma_1\}$  of  $\text{Ca}(x|0, \sigma_2)$  from  $\text{Ca}(x|0, \sigma_1)$  is

$$\int_{-\infty}^{\infty} \text{Ca}(x|0, \sigma_1) \log \frac{\text{Ca}(x|0, \sigma_1)}{\text{Ca}(x|0, \sigma_2)} dx = \log \frac{1}{4\sigma_1\sigma_2} + 2 \log(\sigma_1 + \sigma_2).$$

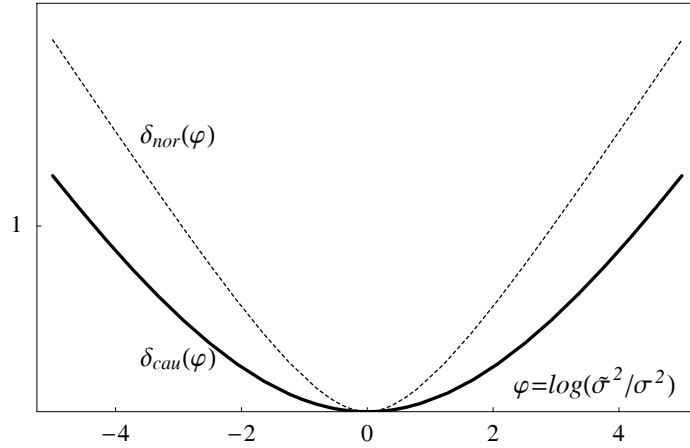


**Figure 4:** Reference posterior distributions of  $\sigma$  for a set of 25 random samples of size  $n = 5$  generated from a  $\text{Ca}(x|0, 2)$  distribution.

Since, in this case,  $\kappa\{\sigma_2 | \sigma_1\} = \kappa\{\sigma_1 | \sigma_2\}$ , the intrinsic discrepancy loss from using  $\tilde{\sigma}$  as a proxy for  $\sigma$  is (Def. 2)

$$\delta\{\tilde{\sigma}, \sigma\} = \log \frac{1}{4\tilde{\sigma}\sigma} + 2 \log(\tilde{\sigma} + \sigma) = \log \frac{1}{4\sqrt{\phi}} + \log(1 + \sqrt{\phi}), \quad (54)$$

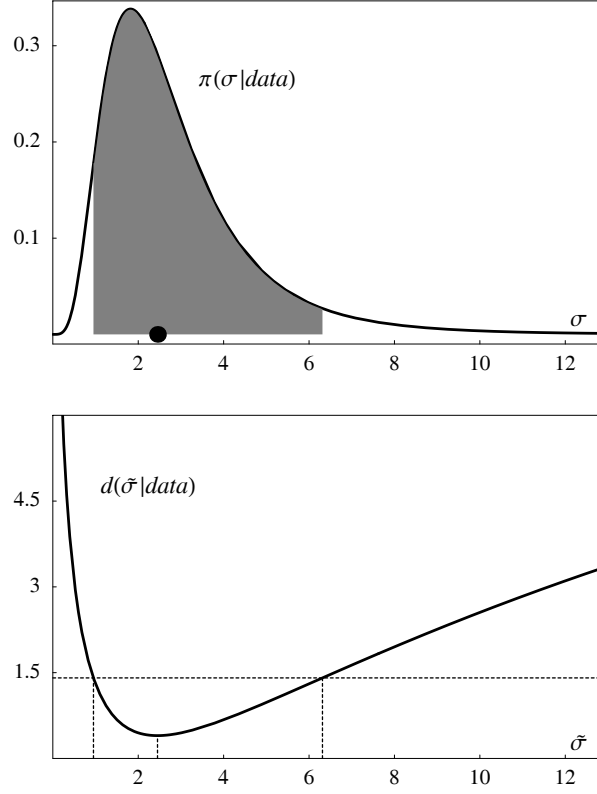
where  $\phi = \tilde{\sigma}^2/\sigma^2$ . Thus, as in the normal case (Eq. 33), the intrinsic discrepancy loss only depends on the variance ratio  $\phi = \tilde{\sigma}^2/\sigma^2$ .



**Figure 5:** Intrinsic discrepancy loss from using  $\tilde{\sigma}$  as a proxy for  $\sigma$  as a function of  $\psi = \log(\tilde{\sigma}^2/\sigma^2)$  for Cauchy (solid line) and normal (dotted line) distributions.

Figure 5 provides a direct comparison between the intrinsic discrepancy loss for the scale parameter in the Cauchy and in the normal case. As one might expect, for any

given value of the ratio  $\phi = \tilde{\sigma}^2/\sigma^2$ , the intrinsic loss is smaller in the Cauchy case than it is in the normal case.



**Figure 6:** Reference posterior density (upper panel) and intrinsic expected loss (lower panel) for the scale parameter  $\sigma$  of a Cauchy  $\text{Ca}(x|0, \sigma)$  distribution, given  $\mathbf{x} = \{-1.78, -0.75, -2.44, -3.30, 8.48\}$ . The intrinsic estimator is  $\tilde{\sigma}_{\text{int}} = 2.452$  (solid dot) and the intrinsic 0.90-credible region is  $C_{0.90}^{\text{int}} = (0.952, 6.314)$  (shaded region).

The intrinsic expected loss from using  $\tilde{\sigma}$  is the reference posterior expectation of the intrinsic discrepancy loss,

$$d(\tilde{\sigma} | \mathbf{x}) = n \int_0^\infty \delta\{\tilde{\sigma}, \sigma\} \pi(\sigma | \mathbf{x}) d\sigma,$$

where  $\delta\{\tilde{\sigma}, \sigma\}$  and  $\pi(\sigma | \mathbf{x})$  are respectively given by (54) and (53), and may easily be computed by numerical integration. The intrinsic estimator of  $\sigma$  is

$$\tilde{\sigma}_{\text{int}}(\mathbf{x}) = \arg \inf_{\tilde{\sigma} > 0} d(\tilde{\sigma} | \mathbf{x})$$

and the  $p$ -credible intrinsic region is the solution  $C_p^{int}(\mathbf{x}) = (\sigma_0, \sigma_1)$  to the equations system

$$\left\{ d(\sigma_0 | \mathbf{x}) = d(\sigma_1 | \mathbf{x}), \quad \int_{\sigma_0}^{\sigma_1} \pi(\sigma | \mathbf{x}) d\sigma = p \right\}.$$

As a numerical illustration, a random sample of size  $n = 5$  was generated from a Cauchy  $\text{Ca}(x|0, 2)$ , yielding  $\mathbf{x} = \{-1.78, -0.75, -2.44, -3.30, 8.48\}$ . The results from the reference analysis of this data set are encapsulated in Figure 6. The reference posterior distribution  $\pi(\sigma | \mathbf{x})$  is represented in the upper panel, and the expected intrinsic loss  $d(\tilde{\sigma} | \mathbf{x})$  from using  $\tilde{\sigma}$  as a proxy for  $\sigma$  is represented in the lower panel. The intrinsic estimator, represented by a solid dot, is  $\tilde{\sigma}_{int} = 2.452$ , and the intrinsic 0.90-credible interval, represented by a shaded region, is  $C_{0.90}^{int} = (0.952, 6.314)$ .

Neither exact Bayesian credible regions nor exact frequentist confidence intervals may be analytically obtained in this problem. The frequentist coverage of the intrinsic credible regions was however analyzed by simulation. A set of 5,000 random samples of size  $n = 5$  were generated from a Cauchy  $\text{Ca}(x|0, 2)$ , and their corresponding intrinsic 0.90-credible intervals were computed; it was then found that the proportion of those intervals which contained the true value  $\sigma = 2$  was 0.905. With 25,000 random samples this proportion was 0.902. This suggests that (as in the normal case) the expected frequentist coverage of reference  $p$ -credible regions, the limit of this algorithm as the number of generated random samples increases, is exactly  $p$ . To further explore this suggestion, a set of 10,000 random samples of size  $n$  were generated from a Cauchy distribution  $\text{Ca}(x|0, \sigma)$  for each of several combinations  $\{n, \sigma\}$  of sample size  $n$  and true value  $\sigma$  of the scale parameter and the corresponding intrinsic  $p$ -credible regions were computed for  $p = 0.90$  and  $p = 0.95$ . Table 2 describes the proportion of these regions which actually contained the value of  $\sigma$  from which the samples had been generated.

**Table 2:** Proportions of intrinsic  $p$ -credible intervals which contained the true value of  $\sigma$  among 10,000 random samples generated from each of several combinations of sample size  $n$  and true value of  $\sigma$ .

	$p = 0.90$			$p = 0.95$		
	$n$			$n$		
$\sigma$	2	12	30	2	12	30
0.5	0.9002	0.9044	0.8999	0.9490	0.9491	0.9507
2.0	0.8971	0.8971	0.9003	0.9467	0.9517	0.9490
4.0	0.9006	0.8960	0.8990	0.9484	0.9497	0.9507

Examination of this table provides strong statistical evidence that the frequentist coverage of reference  $p$ -credible regions is indeed exactly equal to  $p$  for all sample sizes. Indeed, treating each simulation as a Bernoulli trial, the reference posterior

distribution of the frequentist coverage  $\theta_{ij}$  which corresponds to the  $(i, j)$  cell is approximately normal with mean observed proportion  $p_{ij}$  quoted in the table, and standard deviation  $(0.90 * 0.10/10000)^{1/2} = 0.0030$  for the 0.90-credible intervals, and  $(0.95 * 0.05/10000)^{1/2} = 0.0022$  for the 0.95-credible intervals. This makes the respective nominal values 0.90 and 0.95 clearly compatible with the observed results. Notice that this is *not* an asymptotic analysis, as in probability matching theory (Datta and Sweeting, 2005), for it even applies to the smallest possible samples, those with  $n = 2$ .

## References

- Berger, J. O. and Bernardo, J. M. (1992). On the development of reference priors. *Bayesian Statistics 4* (J. M. Bernardo, J. O. Berger, A. P. Dawid and A. F. M. Smith, eds.) Oxford: University Press, 35-60 (with discussion).
- Bernardo, J. M. (1979). Reference posterior distributions for Bayesian inference. *Journal of the Royal Statistical Society*, 41, 113-147 (with discussion). Reprinted in *Bayesian Inference* (N. G. Polson and G. C. Tiao, eds.) Brookfield, VT: Edward Elgar, 1995, 229-263.
- Bernardo, J. M. (2005a). Reference analysis. *Handbook of Statistics 25* (D. K. Dey and C. R. Rao, eds.). Amsterdam: Elsevier, 17-90.
- Bernardo, J. M. (2005b). Intrinsic credible regions: An objective Bayesian approach to interval estimation. *Test*, 14, 317-384 (with discussion).
- Bernardo, J. M. (2006). Intrinsic point estimation of the normal variance. *Bayesian Statistics and its Applications* (S. K. Upadhyay, U. Singh and D. K. Dey, eds.) New Delhi: Anamaya, 110-121.
- Bernardo, J. M. and Juárez, M. (2003). Intrinsic estimation. *Bayesian Statistics*, 7, 465-476.
- Bernardo, J. M. and Rueda, R. (2002). Bayesian hypothesis testing: A reference approach. *International Statistical Review*, 70, 351-372.
- Bernardo, J. M. and Smith, A. F. M. (1994). *Bayesian Theory*. Chichester: Wiley. (2nd edition in preparation).
- Brown, L. (1968). Inadmissibility of the usual estimators of scale parameters in problems with unknown location and scale parameters. *The Annals of Mathematical Statistics*, 39, 24-48.
- Brewster, J. F. and Zidek, J. V. (1974). Improving on equivariant estimators. *Annals of Statistics*, 2, 21-38.
- Brown, L. (1990). Comment on Maata and Casella (1990).
- Datta, G. S. and Sweeting, T. J. (2005). Probability matching priors. *Handbook of Statistics*, 25 (D. K. Dey and C. R. Rao, eds.). Amsterdam: Elsevier, 91-114.
- Fernández, C. and Steel, M. F. J. (1999). Reference priors for the general location-scale model. *Statistics and Probability Letters*, 43, 377-384.
- Juárez, M. A. (2004). *Métodos Bayesianos Objetivos de Estimación y Contraste de Hipótesis*. Ph.D. Thesis, Universidad de Valencia, Spain.
- Kullback, S. and Leibler, R. A. (1951). On information and sufficiency. *The Annals of Mathematical Statistics*, 22, 79-86.
- Maata, J. M. and Casella, G. (1990). Developments in decision-theoretic variance estimation. *Statistical Science*, 5, 90-120, (with discussion).
- Robert, C. P. (1996). Intrinsic loss functions. *Theory and Decision*, 40, 192-214.
- Rukhin, A. L. (1987). How much better are better estimators of a normal variance. *Journal of The American Statistical Association*, 82, 925-928.

Schervish, M. J. (1995). *Theory of Statistics*. Berlin: Springer

Stein, C. (1964). Inadmissibility of the usual estimator for the variance of a normal distribution with unknown mean. *Annals of the Institute of Statistical Mathematics*, 16, 155-160.

**Discussion of “Objective bayesian  
point and region estimation in  
location-scale models”  
by José M. Bernardo**





**Miguel A. Gómez Villegas**

Departamento de Estadística e Investigación Operativa

Universidad Complutense de Madrid

Let me begin by congratulating Professor Bernardo for his excellent job in objective Bayesian analysis. This paper, and the closely related Bernardo (2005), present a unified theory of estimation by point and credible regions based on information ideas he has used previously to define reference priors. The idea originates from the study of both problems as decision problems, where the loss function is the “intrinsic discrepancy” inspired in the Kullback-Leibler divergence, and defined as the minimum of  $k_x\{\tilde{\theta}, \tilde{\lambda}|\theta, \lambda\}$  and  $k_x\{\theta, \lambda|\tilde{\theta}, \tilde{\lambda}\}$  where

$$k_x\{\tilde{\theta}, \tilde{\lambda}|\theta, \lambda\} = \int_{\chi(\theta, \lambda)} \pi(x|\theta, \lambda) \ln \frac{\pi(x|\theta, \lambda)}{\pi(x|\tilde{\theta}, \tilde{\lambda})} dx$$

An intrinsic point estimator is then defined as the Bayes estimator which corresponds to the intrinsic loss and the appropriate reference prior. A  $p$ -credible intrinsic region estimator is defined as the lowest posterior loss  $p$ -credible with respect to the intrinsic loss and the appropriate reference prior.

A first question is: do we need to employ

$$\int_{C_p^{int}} \pi(\theta|x) d\theta \geq p$$

with the inequality instead of equality to allow the discrete case?

Second, it would be useful to have a better understanding of the proposed approach to applying these ideas to the exponential distribution family instead of location-scale models; this is a family of distributions greater than the other.

Professor Bernardo claims that in one-dimensional problems, one may define probability centred credible intervals, and these are invariant under reparametrization. Will it not be necessary to suppose that the transformation is monotonic?

Third, on a more philosophical basis, I think that invariance is a compelling argument for point estimations and for credible regions. Indeed both point estimations and credible regions are two answers to the same question: how we can eliminate the uncertainty

about  $\theta$ . Bernardo's approach permits one to obtain invariance under reparametrization in both problems.

Fourth, the chosen examples show the coherence between frequentist inference and Bayesian inference. When intrinsic credible regions that require minimal subjective inputs are employed, exact frequentist confidence regions are obtained, at least in the normal mean and variance. This fact is similar to the one obtained by this discussant in Gómez-Villegas and González-Pérez (2005) and references therein. I wonder if Professor Bernardo has any idea about the essential reasons behind the matching properties between intrinsic credible regions and confidence regions in these cases?

Fifth, adopting this approach to credible set construction, I see problems in computations, the posterior intrinsic loss integrated over a large dimensional space. From the point of view of applications, a simple asymptotic approximation to normality should be necessary.

In closing, I would like to thank the editor of the journal for giving me the opportunity to discuss this paper.

## References

- Bernardo, J. M. (2005). Intrinsic credible regions: an objective Bayesian approach to interval estimation. *Test*, 14, 317-384.
- Gómez-Villegas, M. A. and González-Pérez, B. (2005). Bayesian analysis of contingency tables. *Communications in Statistics-Theory and Methods*, 34, 1743-1754.

**Dennis V. Lindley**

ThomBayes@aol.com

Two concepts are basic to the ideas of this excellent paper: *objectivity* and the concept of estimation as a *decision* problem. In the author's skilful hands, these lead to reference priors and intrinsic loss functions, and hence, by minimizing expected loss, to estimates which are often superior to the conventional ones. It can be said with some confidence that we have here a solution to the problem Harold Jeffreys first posed around 1939 of providing an objective, coherent method for scientific inference. The development employs several subtle ideas, and considerable mathematical complexity, but one feature that struck me is that the final results are usually fairly simple and look right. An example of this is provided by the loss functions in Figure 3, which have the reasonable convexity property around the true value but, unlike quadratic loss, exhibit sensible concavity at more discrepant values. I would have preferred the loss to have been bounded but, with normal distributions and their thin tails, this scarcely matters. To be bounded may be more important with the fat tails of the Cauchy in Figure 6, in order to avoid paradoxes of the St. Petersburg type. A related point is that although the mathematics can be formidable, at least in the view of some applied statisticians, once it has been done the practitioner can easily use the results in the confidence that the machinery used to produce them is sound. It is comparable to driving a car, without knowing how it was made, but having confidence in the manufacturer.

Granted the basic concepts, this is an important paper, but was Jeffreys right to search for objectivity, and was Fisher wrong in dismissing decision concepts from inference? I think Jeffreys was wrong and Fisher was right. At the risk of repeating what I have said before, it seems to me that inference and decision-making are distinct and both are subjective. In other words, the two basic concepts, that provide the foundations of this paper, are suspect.

Consider first the fixed likelihood upon which all the arguments in the paper rest. Is it really objective? There are a few cases where substantial evidence for normality exists, but often the normal, or another member of the exponential family, is used merely for mathematical simplicity. With the increased computing power available today, statisticians are less constrained and can use other distributions that appear more realistic, thereby introducing subjectivity. There are some popular data sets that have been repeatedly analysed using different likelihoods. Where is the objectivity there? It is interesting that Bernardo uses one symbol,  $p$ , for probabilities of data but another,

$\pi$ , for probabilities of parameters. In reality both  $p$  and  $\pi$  reflect beliefs about data and parameters respectively, obey the same rules and do not deserve separate treatments.

Inference concerns parameters. (It is more practical to make inference about future data but I do not explore that trail here.) What are these parameters, the  $\theta$  and  $\lambda$  of the paper? If our statistical analyses are to be of use in data analysis,  $\theta$  at least ought to relate to something in the real world. Bernardo has only one sentence about this, referring to  $\theta$  as the age of the earth. Putting aside intelligent designers, reputable scientists differ in their views of the age. In other words, their ideas are subjective, so that before relevant data about the age are considered, their different views need to be included. Another relevant fact is that information about the age of the earth does not come from data with normal, or any other objective, likelihood. More conspicuous examples of subjectivity are apparent with clinical trials, where the different views of drug companies and official bodies are consulted before the trial. This became clear recently when a trial went horribly wrong and experts claimed the probabilities used were, in their opinion, unsound. So often today,  $\theta$  is regarded as nothing more than a symbol, whereas, to be of value, it has to refer to reality and hence influenced by opinions about that reality. These opinions should be incorporated into the analysis, not ignored and replaced by a reference prior, especially when this is improper.

All of us have, at some time, expressed an opinion about something without having any intention of basing any action upon that opinion. In statistics, this opinion-forming is inference and means we infer the value of the real  $\theta$ . In the Bayesian paradigm this is done by means of your probability distribution of  $\theta$ , given the data and the original information about  $\theta$ . Whilst it is true that any inference has to be capable of being used as a basis for action, for otherwise what use is it, it is not true that inference has to have immediate actions in mind. In particular, inference does not require a loss function, and certainly not a loss function that ignores reality. In Bayesian terms, there is only one inference, the posterior distribution and, although it may be advantageous to summarize its main features, such approximations scarcely need elaborate techniques, except in the case of many parameters.

Inference from data consists in modelling that data in the form of a likelihood depending on parameters, supplying your opinion of the parameters prior to the data, and combining likelihood and prior by Bayes theorem. Finally the nuisance aspect of the parameters is removed by integration. When several people are involved there may be disagreements over likelihood or prior. These may be removed by discussion but, if this fails, the calculations may be repeated under different subjective opinions and the posteriors compared. That science is objective is a myth. Apparent objectivity in science only arises when the data are extensive.

This paper explores a field that, in my view, is not in the broad stream of statistics. This is not to deny it great merit, for we now know what that field contains, material of real merit from which all can learn.

**Mark J. Schervish**

Carnegie Mellon University, USA

I admire Professor Bernardo for his steadfastness and resolution in staying the course of research into reference priors and other so-called objective Bayesian methods. Despite repeated attacks dating back to the discussion of Bernardo (1979) he has continually risen to the challenge of making these methods palatable to practitioners and theoreticians alike. I will not here rehearse all of the criticisms or the support for his work in this area. I refer the interested reader to the various discussions of the papers listed in the reference list to Professor Bernardo's paper. I will mention just a few problems that I have with the methods as well as what I like about them.

To begin with a positive note, I like the idea of having a transformation-equivariant estimation procedure for non-decision-theoretic inference. When one is faced with a decision problem in which a specific loss function is relevant, then one does not care whether one's inference satisfies an ad hoc criterion such as transformation equivariance. On the other hand, when one merely wishes to report an estimate of some quantity, especially the parameter of a statistical model which most likely is a figment of one's imagination (model) anyway, then it becomes difficult to explain why the estimate of an equivalent parameter is not the equivalent estimate. Indeed, I believe that the intrinsic discrepancy loss satisfies a slightly stronger invariance than is stated in (10). I believe that one could apply a one-to-one reparameterization of the form  $\phi = \phi(\theta)$  and  $\psi = \psi(\lambda, \theta)$  and still achieve (10). Of course, a completely general reparameterization would change the meaning of the parameter of interest, and yet the desire for an equivariant estimate would remain.

One of the serious concerns with reference priors is their violation of the likelihood principle. The reference priors are different for binomial sampling and negative binomial sampling so that even if the observed data could have come from either sampling scheme, the posterior would depend on the sampling plan. If one were to observe a binomial sample and use the reference prior, and later observe a negative binomial sample, one would get a different inference than if one were to observe the same two samples in the other order. As mentioned earlier, various discussants have described other concerns with the methods advocated in the manuscript, and I will let the reader find them in their original forms. I will add only one other concern that I have, and that is with the use of the description of these methods as "objective". I suppose that, so long as one agrees with all of the reasons put forth for why such methods should be

used, then one will use the methods and they become objective in that sense. But any set of methods could be called objective on those grounds. One of the main strengths of Bayesian methodology is that it forces users to be explicit about the assumptions that they are making. People who think that they are using objective methods are simply borrowing a collection of subjective assumptions and ignoring the fact that choices were made by someone else arriving at those assumptions. When you lay your assumptions out for all to see, you are in a position to evaluate the sensitivity of your inferences to the assumptions. If you hide behind a cloak of objectivity, you may produce the same answer that others produce, but you have lost the ability to see what is the effect of the subjective choices that were made.

## Rejoinder

I am extremely grateful to the three discussants by their thoughtful comments. I will answer them individually.

*Gómez-Villegas.* If the parameter of interest  $\theta$  is discrete, then we would certainly need to work with regions  $C$  such that  $\int_C \pi(\theta|\theta) d\theta \geq p$  since, in that special case, not all credible probabilities  $p$  would be attainable. However, point and region estimation are usually done with *continuous* parameter spaces, and this is indeed the case in the location and scale models considered in this paper. In that situation, the equality may always be obtained.

The ideas discussed in the paper may certainly be applied to models in the (generalized) exponential family and it is likely that this would lead to some rather general results. I did not have time and space to do this here, but it is certainly a research line well worth exploring.

As Professor Gómez-Villegas points out, the invariance arguments invoked only refer to monotonic, one-to-one transformations of the parameter. Even though not always explicitly stated, we were indeed always assuming this to be the case.

I believe that the *exact* numerical coincidence between objective credible regions and frequentist confidence interval is the exception, not the rule; when it happens, it is the consequence of the existence of pivotal quantities, so that the reference distribution of the pivot (considered as a function of the parameter) is precisely the same as its sampling distribution (considered as a function of the data). In particular, this coincidence cannot exist if data are discrete, as in the case of binomial or Poisson data. Beyond the particular situations where pivots exist, one may only expect an asymptotic approximation: objective credible regions are typically *approximate* confidence intervals, the approximation improving with the sample size.

Routine application of the methods described in this paper will certainly require either available software producing the exact results (not difficult to write in the standard examples which constitute the vast majority of applications) and/or appropriate analytical approximations. The latter may easily be obtained, as in the examples contained in the paper, by using the normal approximation with the parametrization induced by the appropriate variance-stabilizing transformation, and then making use of the invariance properties of the procedures.

*Lindley.* I am really proud that Professor Lindley may believe that the procedures described provide an objective coherent method for scientific inference in the sense demanded by Jeffreys, and I am very grateful for that comment.

It would certainly be better from a foundations viewpoint if the expected loss were bounded, but information measures with continuous parameters are not bounded (one needs infinite amount of information to know precisely a real number) and yet have all kind of attractive properties.

To repeat in print the basics of an argument that Professor Lindley and I have often had in private conversations,

- (i) I believe, with Jeffreys, that Fisher was wrong in dismissing decision concepts in inference. If, by some reason, you must choose an estimate, then (whether you like it or not) you have a well posed decision problem where the action space is the set of parameter values; then foundations dictate that (to act rationally) you *must* use a loss function. For instance, in one continuous parameter problems, the median may well be an estimate with good robustness properties, but the fact remains that this would be a good estimate if (*and only if*) your loss function is well approximated by a linear, symmetric loss function.
- (ii) I applaud the use of subjective priors when the problem is simple and small enough for the required probability assessments to be feasible (which is *not* frequent). But, even in this case, there is no reason while other people should necessarily accept a subjective prior which goes beyond clearly stated assumptions and verifiable (possibly historical) data. There is a clear need for some commonly accepted minimum set of conclusions to be solely derived from assumptions and data, and this is precisely what reference posteriors provide. As their name indicate, they are proposed as a *reference*, to be compared with subjective posteriors when these are available. This is part of a necessary exercise in sensitivity analysis, by making explicit which parts of the conclusions depend on a particular subjective prior, and which parts are implied by the model assumed and the data obtained.

As Professor Lindley points out, although inferential statements are typically used as a basis for action, there are many situations where inferences are to be drawn without any specific action in mind. This is precisely why we suggest the use of an information-based loss function. If a particular action is in mind, one should certainly use a context dependent loss function which appropriately describes the decision problem analyzed. If no particular decision problem is in mind, one is bound to use some conventional loss function. We have argued that conventional loss functions (such as the ubiquitous quadratic loss) are often unsatisfactory. Instead, for “pure inference” problems one should try to minimize the information loss due to the use of an estimate of the unknown parameter value; and this, I believe, is appropriately captured by the intrinsic discrepancy loss.

*Schervish.* I am very glad to read that Professor Schervish appreciates the importance of invariant procedures. In teaching, I often start my lectures by stating that any inferential



procedure which is not invariant under monotonic transformations of the parameter is suspect, and go on to provide a set of examples of those as “counterexamples” to common statistical procedures.

I agree with Professor Schervish on the importance of the likelihood principle, but I believe that the principle is actually compatible with a sensible use of reference distributions. Indeed, a reference posterior encapsulates, by definition, the (minimal) inferential statements you could proclaim about the parameter of a model *if* your prior was that maximizing the information that data generated from that *particular* model could possibly provide. If you change the model (even if the new model induces a proportional likelihood function), you change the reference prior. Thus, different reference posteriors corresponding to different sampling schemes with Bernoulli observations provide a collection of *conditional* answers (one for each sampling scheme one is willing to consider), which may all be part of the sensitivity analysis to changes in the prior mentioned above.

Objectivity is indeed an emotionally charged word, and it should be explicitly qualified whenever it is used. No statistical analysis is seriously objective, if only because the choice of both the experiment design and the model used have typically very strong subjective inputs. However, the frequentist paradigm is sold as “objective” just because its conclusions are only conditional on the model assumed and the data obtained, and this objectivity illusion has historically helped frequentist to keep a large share of the statistics market. I claim for the procedures described in this paper the right to use “objective” in precisely the same sense: these are procedures which are only conditional on the assumed model and the observed data. The use of the word “objective” in this precise, limited sense may benefit, I believe, the propagation of the Bayesian paradigm. For a recent discussion of this and related issues see Berger (2006) and ensuing discussion.

I fully agree with Professor Schervish on the paramount importance of clearly presenting the assumptions needed for an inferential statement. In the case of reference posteriors this should typically read as a *conditional* statement of the form: “*If* available data  $x$  had been generated by model  $\mathcal{M} \equiv \{p_x(\cdot|\omega), \omega \in \Omega\}$  and prior information about  $\theta(\omega)$  were minimal with respect to the information about  $\theta(\omega)$  that repeated sampling from  $\mathcal{M}$  could possibly provide *then*, the marginal reference posterior  $\pi(\theta|x)$  encapsulates what could be said about the value of  $\theta$ , solely on the basis of that information”.

## References

- Berger, J. O. (2006). The case for objective Bayesian analysis. *Bayesian Analysis*, 1, 385-402 and 457-464 (with discussion).

